

# On the Total Variation Distance of Labelled Markov Chains

Taolue Chen

Middlesex University London

t.l.chen@mdx.ac.uk

Stefan Kiefer

University of Oxford

stekie@cs.ox.ac.uk

## Abstract

Labelled Markov chains (LMCs) are widely used in probabilistic verification, speech recognition, computational biology, and many other fields. Checking two LMCs for equivalence is a classical problem subject to extensive studies, while the total variation distance provides a natural measure for the “inequivalence” of two LMCs: it is the maximum difference between probabilities that the LMCs assign to the same event.

In this paper we develop a theory of the total variation distance between two LMCs, with emphasis on the algorithmic aspects: (1) we provide a polynomial-time algorithm for determining whether two LMCs have distance 1, i.e., whether they can almost always be distinguished; (2) we provide an algorithm for approximating the distance with arbitrary precision; and (3) we show that the threshold problem, i.e., whether the distance exceeds a given threshold, is NP-hard and hard for the square-root-sum problem. We also make a connection between the total variation distance and *Bernoulli convolutions*.

**Categories and Subject Descriptors** G.3 [Probability and Statistics]; D.2.4 [Software/Program Verification]

**General Terms** Theory

**Keywords** Labelled Markov Chains, Total Variation Distance

## 1. Introduction

A (discrete-time, finite-state) *labelled Markov chain (LMC)* has a finite set  $Q$  of states and for each state a probability distribution over its outgoing transitions. Each outgoing transition is labelled with a letter from a given finite alphabet  $\Sigma$ , and leads to a target state. Figure 1 depicts two LMCs. The semantics is as follows: The chain starts in a given initial state (or in a random state according to a given initial distribution), picks a random transition according to the state’s distribution over the outgoing transitions, outputs the letter of the transition, moves to the target state, and repeats. In such a way, the chain produces a random infinite sequence of letters, i.e., a random infinite word. We regard this infinite word as “observable” to the environment, whereas the infinite sequence of states remains “internal” to the chain. Formally, an LMC defines a probability space whose samples are infinite words (also called *runs* later) over  $\Sigma$ . In [19], it is classified as a generative model. LMCs

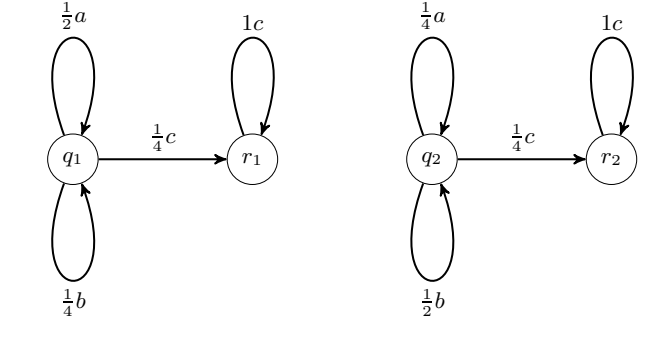


Figure 1. Two LMCs.

appear as *hidden Markov models* in speech recognition and in several areas of computational biology, cf. [14]. LMCs, sometimes in the form of *probabilistic automata* [17], are also fundamental for modelling probabilistic systems.

Checking whether two LMCs (or, similarly, two probabilistic automata) are (*language*) *equivalent* is a classical problem, going back to the seminal work of Schützenberger [18] and Paz [15]. More recently, this problem was revisited, as various verification problems on probabilistic systems can be reduced to it (see, e.g., [12]). As a consequence, efficient polynomial-time algorithms and tools for equivalence checking have been developed [4, 6, 11, 12]. If two systems are found to be *not* equivalent, the question arises on *how different* they are. The *distance* of two LMCs provides a measure for their difference, with the extreme cases being distance 0 for equivalence and distance 1 for (almost-sure) distinguishability. The *total variation distance*, which is a standard distance measure [9] between two probability distributions, yields a natural measure of the distance of two LMCs. Given two probability distributions  $\pi_1$  and  $\pi_2$  over the same *countable* set  $\Omega$ , the total variation distance is defined as

$$d(\pi_1, \pi_2) := \max_{E \subseteq \Omega} |\pi_1(E) - \pi_2(E)|. \quad (1)$$

In words,  $d(\pi_1, \pi_2)$  is the largest possible difference between probabilities that  $\pi_1$  and  $\pi_2$  assign to the same event. Furthermore, we have  $d(\pi_1, \pi_2) = \pi_1(E) - \pi_2(E)$  for

$$E = \{r \in \Omega \mid \pi_1(r) \geq \pi_2(r)\}, \quad (2)$$

so the event  $E$  is a maximizer in (1). The total variation distance is—up to a factor of 2—equal to the  $L_1$ -norm of the difference between  $\pi_1$  and  $\pi_2$ :

$$2d(\pi_1, \pi_2) = \|\pi_1 - \pi_2\|_1 := \sum_{x \in \Omega} |\pi_1(x) - \pi_2(x)|.$$

When applying the total variance distance to LMCs, it should be emphasized that the sample space  $\Omega = \Sigma^\omega$  (i.e., the set of infinite words over  $\Sigma$ ) is *uncountable*. Hence the maximum in

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.

CSL-LICS 2014, July 14–18, 2014, Vienna, Austria.  
Copyright © 2014 ACM 978-1-4503-2886-9...\$15.00.  
<http://dx.doi.org/10.1145/2603088.2603099>

the definition of (total variation) distance needs to be replaced by the supremum. Concretely, assume two LMCs  $\mathcal{M}_1, \mathcal{M}_2$  with initial state distributions, the LMCs assign each (measurable) event  $E \subseteq \Sigma^\omega$  a probability  $\pi_1(E)$  and  $\pi_2(E)$ , respectively. So the (total variation) distance between  $\mathcal{M}_1, \mathcal{M}_2$  is defined as

$$d(\pi_1, \pi_2) := \sup_{E \subseteq \Sigma^\omega} |\pi_1(E) - \pi_2(E)|.$$

It is not clear a priori if a maximizer event exists. We will show later in this paper that it does exist. In particular this means that  $d(\pi_1, \pi_2) = 1$  holds if and only if there is an event  $E$  with  $\pi_1(E) = 1$  and  $\pi_2(E) = 0$ .

While being an intriguing theoretical question, the study of the distance between LMCs also has practical implications. For instance, in the verification of anonymity properties [11, 12] the following scenario is common: Two users are modelled as LMCs and leave a trace (i.e., emit a run). An evil agent knows the two users, and sees a single trace. The agent wants to find out which of the two users has emitted the trace. Clearly language equivalence (distance 0) of LMCs implies anonymity of the users. If the distance is nonzero, one may ask if the agent can identify the users almost surely. If the distance is 1, the agent succeeds with probability 1, because the agent can define an event  $E$  that occurs in the first LMC with probability 1, and in the second one with probability 0; all the agent has to do is to check whether the given run belongs to  $E$ . Conversely, if the distance is less than 1, the agent cannot almost-surely distinguish the users. From this point of view, a distance less than 1 is a minimum requirement for some form of user anonymity, which could perhaps be called *deniability*.

Another example is probabilistic model checking where computing the probability of certain events  $E$  is of central interest. If the distance between some given LMCs is small (and known or bounded above), computing the probability of  $E$  in one of those chains may be enough for obtaining good bounds on the probability of  $E$  in the other chains. This may lead to savings in the overall model-checking time.

**Main Contributions.** In this paper we develop a theory for the total variation distance between two LMCs. We pay special attention to the algorithmic and computational aspects of the problem. We make the following contributions:

- (1) We demonstrate some basic properties of the total variation distance between two LMCs: (a) the supremum in the definition can be “achieved”, and we exhibit a maximizing event, although we show that the maximizing event is *not*  $\omega$ -regular in general; (b) the distance of two LMCs can be irrational even if all probabilities appearing in their description are rational.
- (2) We study the *qualitative* variant of the distance problem, i.e., to decide whether two LMCs have distance 1 or 0. The distance-0 problem amounts to the language equivalence problem for probabilistic automata, for which a polynomial-time algorithm exists. We provide a polynomial-time algorithm for the distance-1 problem.
- (3) We study the *quantitative* variant of the distance problem. In light of (1), at best one can hope to *approximate* the distance rather than to really *compute* it (at least in the classical complexity theory framework). To this end, we provide an algorithm for approximating the distance with arbitrary precision. We also link the problem to *Bernoulli convolutions* by providing an LMC where the distance of two states of this LMC is related to *Bernoulli convolutions*, thus indicating the intricacy of the distance.
- (4) We study the threshold problem, i.e., to decide whether the distance exceeds a given threshold. While leaving decidability

of the problem open, we show that the problem is both NP-hard and hard for the square-root-sum problem.

**Structure of the Paper.** In Section 2 we provide technical preliminaries. In Section 3 we give two examples for LMCs and their distances. In Section 4 we discuss two sequences that converge to the distance from below and from above, yielding an approximation algorithm. In Section 5 we show that an event with maximum difference in probabilities always exists, and we exhibit such a “witness” event. In Section 6 we show that the distance can be irrational, and we give lower complexity bounds for the threshold problem. In particular, in Section 6.1 we exhibit an LMC where the distance depends on the probabilities in the LMC in intricate ways, as witnessed by a connection to *Bernoulli convolutions*. In Section 7 we develop a polynomial-time algorithm for deciding whether two LMCs have distance 1. In Section 8 we discuss related work. Finally, in Section 9 we offer some conclusions and highlight open problems. Missing proofs can be found in an appendix.

## 2. Preliminaries

We write  $\mathbb{N}$  for the set of nonnegative integers.

Let  $Q$  be a finite set. By default we view *vectors*, i.e., elements of  $\mathbb{R}^Q$ , as row vectors. For a vector  $\mu \in [0, 1]^Q$  we write  $|\mu| := \sum_{q \in Q} \mu(q)$  for its  $L_1$ -norm. A vector  $\mu \in [0, 1]^Q$  is a *distribution* (resp. *subdistribution*) over  $Q$  if  $|\mu| = 1$  (resp.  $|\mu| \leq 1$ ). For  $q \in Q$  we write  $\delta_q$  for the (*Dirac*) distribution over  $Q$  with  $\delta_q(q) = 1$  and  $\delta_q(r) = 0$  for  $r \in Q \setminus \{q\}$ . For a subdistribution  $\mu$  we write  $\text{supp}(\mu) = \{q \in Q \mid \mu(q) > 0\}$  for its support. Given two vectors  $\mu_1, \mu_2 \in [0, 1]^Q$  we write  $\mu_1 \leq \mu_2$  to say that  $\mu_1(q) \leq \mu_2(q)$  holds for all  $q \in Q$ . We view elements of  $\mathbb{R}^{Q \times Q}$  as *matrices*. A matrix  $M \in [0, 1]^{Q \times Q}$  is called *stochastic* if each row sums up to one, i.e., for all  $q \in Q$  we have  $\sum_{r \in Q} M(q, r) = 1$ .

**Definition 1.** A labelled (discrete-time, finite-state) Markov chain (LMC) is a tuple  $\mathcal{M} = (Q, \Sigma, M)$  where

- $Q$  is a finite set of states,
- $\Sigma$  is a finite alphabet of labels, and
- $M : \Sigma \rightarrow [0, 1]^{Q \times Q}$  specifies the transitions, so that  $\sum_{a \in \Sigma} M(a)$  is a stochastic matrix.

Intuitively, if the LMC is in state  $q$ , then with probability  $M(a)(q, q')$  it emits  $a$  and moves to state  $q'$ . For the complexity results of this paper, we assume that all the numbers in the matrices  $M(a)$  for  $a \in \Sigma$  are rationals given as fractions of integers represented in binary. We extend  $M$  to the mapping  $M : \Sigma^* \rightarrow [0, 1]^{Q \times Q}$  with  $M(a_1 \cdots a_k) = M(a_1) \cdots M(a_k)$  for  $a_1, \dots, a_k \in \Sigma$ . Intuitively, if the LMC is in state  $q$  then with probability  $M(w)(q, q')$  it emits the word  $w$  and moves (in  $|w|$  steps) to state  $q'$ .

Fix an LMC  $\mathcal{M} = (Q, \Sigma, M)$  for the rest of this section. A run of  $\mathcal{M}$  is an infinite sequence  $a_1 a_2 \cdots$  with  $a_i \in \Sigma$  for all  $i \in \mathbb{N}$ . We write  $\Sigma^\omega$  for the set of runs. For a run  $r = a_1 a_2 \cdots$  and  $i \in \mathbb{N}$  we write  $r_i := a_1 a_2 \cdots a_i$ . For a set  $W \subseteq \Sigma^*$  of finite words, we define  $W\Sigma^\omega := \{wu \mid w \in W, u \in \Sigma^\omega\} \subseteq \Sigma^\omega$ ; i.e., the set of runs that have a prefix in  $W$ . For  $w \in \Sigma^*$  we define  $\text{Run}(w) := \{w\}\Sigma^\omega$ ; i.e.,  $\text{Run}(w)$  is the set of runs starting with  $w$ .

To an (initial) distribution  $\pi$  over  $Q$  we associate the probability space  $(\Sigma^\omega, \mathcal{F}, \text{Pr}_\pi)$ , where  $\mathcal{F}$  is the  $\sigma$ -field generated by all *basic cylinders*  $\text{Run}(w)$  with  $w \in \Sigma^*$ , and  $\text{Pr}_\pi : \mathcal{F} \rightarrow [0, 1]$  is the unique probability measure such that  $\text{Pr}_\pi(\text{Run}(w)) = |\pi M(w)|$ . We generalize the definition of  $\text{Pr}_\pi$  to subdistributions  $\pi$  in the obvious way, yielding sub-probability measures. An *event* is a measurable set  $E \subseteq \Sigma^\omega$ . In this paper we consider only measurable subsets of  $\Sigma^\omega$ , and when we write  $E \subseteq \Sigma^\omega$ , the set  $E$  is meant to be measurable. An event is  $\omega$ -regular, if it is equal to a language ac-

cepted by a nondeterministic Büchi automaton. When confusion is unlikely, we may identify the (sub-)distribution  $\pi$  with the induced (sub-)probability measure  $\Pr_\pi$ ; i.e., for events  $E \subseteq \Sigma^\omega$  we may write  $\pi(E)$  for  $\Pr_\pi(E)$ . For a distribution  $\pi$  and a word  $w \in \Sigma^*$ , we write  $\pi^w$  as a shorthand for  $\pi M(w)$ ; intuitively this is the state subdistribution after emitting  $w$ . We have  $\Pr_\pi(\text{Run}(w)) = |\pi^w|$ .

We reserve  $\pi, \rho$  (and  $\pi_1, \pi_2, \dots$ ) for distributions over  $Q$ , often viewing  $\pi_1, \pi_2$  as given initial distributions. Similarly, we reserve  $\mu, \nu$  for subdistributions over  $Q$ . (But note that  $\pi^w$  for  $w \in \Sigma^*$  is a subdistribution in general.)

Given two initial distributions  $\pi_1, \pi_2$ , we define the (*total variation*) distance between  $\pi_1$  and  $\pi_2$  by

$$d(\pi_1, \pi_2) := \sup_{E \subseteq \Sigma^\omega} |\pi_1(E) - \pi_2(E)|.$$

Recall that  $E \subseteq \Sigma^\omega$  implicitly means that  $E$  is measurable. As  $\pi_1(E) - \pi_2(E) = -(\pi_1(\Sigma^\omega \setminus E) - \pi_2(\Sigma^\omega \setminus E))$ , we have in fact  $d(\pi_1, \pi_2) = \sup_{E \subseteq \Sigma^\omega} (\pi_1(E) - \pi_2(E))$ .

**Remark 2.** One could analogously define the total variation distance between two LMCs  $\mathcal{M}_1 = (Q_1, \Sigma, M_1)$  and  $\mathcal{M}_2 = (Q_2, \Sigma, M_2)$  with initial distributions  $\pi_1$  and  $\pi_2$  over  $Q_1$  and  $Q_2$ , respectively. Our definition is without loss of generality, as one can take the LMC  $\mathcal{M} = (Q, \Sigma, M)$  where  $Q$  is the disjoint union of  $Q_1$  and  $Q_2$ , and  $M$  is defined using  $M_1$  and  $M_2$  in the straightforward manner.

We write  $\mu_1 \equiv \mu_2$  to denote that  $\mu_1$  and  $\mu_2$  are (*language*) equivalent, i.e., that  $|\mu_1^w| = |\mu_2^w|$  holds for all  $w \in \Sigma^*$ . The following proposition states in particular that equivalence can be decided in polynomial time, and that equivalence and the distance being zero are equivalent.

**Proposition 3.**

- (a) We have  $\pi_1 \equiv \pi_2$  if and only if  $d(\pi_1, \pi_2) = 0$ .
- (b) One can compute in polynomial time a set  $\mathcal{B} \subseteq \mathbb{Q}^{2|Q|}$  of column vectors, with  $|\mathcal{B}| \leq 2|Q|$ , such that for all subdistributions  $\mu_1, \mu_2$  we have  $\mu_1 \equiv \mu_2$  if and only if  $(\mu_1 \mu_2) \cdot b = 0$  holds for all  $b \in \mathcal{B}$ . Here,  $(\mu_1 \mu_2) \in [0, 1]^{2|Q|}$  is the row vector obtained by gluing  $\mu_1, \mu_2$  together. (Note that  $(\mu_1 \mu_2) \cdot b$  is a scalar.)
- (c) We have  $\mu_1 \equiv \mu_2$  if and only if  $|\mu_1^w| = |\mu_2^w|$  holds for all  $w \in \Sigma^*$  with  $|w| = 2|Q|$ .
- (d) It is decidable in polynomial time whether  $\mu_1 \equiv \mu_2$  holds. Hence it is also decidable in polynomial time whether  $d(\pi_1, \pi_2) = 0$  holds.

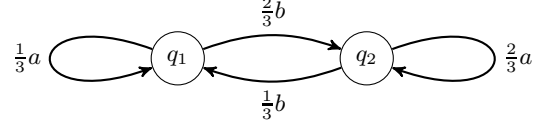
Proposition 3 (a) is immediate from the definitions. Parts (b)–(d) follow from a linear-algebra argument described, e.g., in [6, 15, 18]. We sketch this argument in Appendix A.

### 3. Examples

We illustrate some phenomena of the distance by two examples. The main observations are that the distance of two LMCs can be irrational (Example 1), and in general, they must be differentiated by events which are *not*  $\omega$ -regular (Example 1), even if their distance is 1 (Example 2).

#### 3.1 Example 1

Consider the LMCs from Figure 1 on page 1. As discussed in Remark 2, we can equivalently view them as a single LMC. To illustrate the definitions we study the distance between states  $q_1$  and  $q_2$ , or more precisely, between the Dirac distributions  $\delta_{q_1}$  and  $\delta_{q_2}$ . Note that we have  $\delta_{r_1} \equiv \delta_{r_2}$ , as both  $r_1$  and  $r_2$  keep emitting the letter  $c$ . On the other hand we have  $\delta_{q_1} \not\equiv \delta_{q_2}$  and so  $d(\delta_{q_1}, \delta_{q_2}) > 0$ . With probability 1, one of the states  $r_1, r_2$  will eventually be reached. So events are characterized by the words



**Figure 2.** A 2-state LMC. The two states have distance 1.

over  $a, b$  emitted before the infinite  $c$ -sequence. More formally, for any event  $E \subseteq \Sigma^\omega = \{a, b, c\}^\omega$  one can define  $W_E := \{w \in \{a, b\}^* \mid wc^\omega \in E\}$  so that we have

$$\delta_{q_i}(E) = \sum_{w \in W_E} \delta_{q_i}(\{w\}\{c\}^\omega) \quad \text{for } i \in \{1, 2\}.$$

It is easy to see that  $\delta_{q_1}(\{a\}\{c\}^\omega) = \frac{1}{8}$  and  $\delta_{q_2}(\{a\}\{c\}^\omega) = \frac{1}{16}$ . Consider any event  $E$  with  $W_E$  defined as above. If  $a \in W_E$ , then  $\delta_{q_2}(E) \geq \delta_{q_2}(\{a\}\{c\}^\omega) = \frac{1}{16}$ . If  $a \notin W_E$ , then  $\delta_{q_1}(E) \leq 1 - \delta_{q_1}(\{a\}\{c\}^\omega) = \frac{7}{8}$ . So for any  $E$  we have  $\delta_{q_1}(E) - \delta_{q_2}(E) \leq \max\{1 - \frac{1}{16}, \frac{7}{8}\} = \frac{15}{16}$ . By symmetry we also have  $\delta_{q_2}(E) - \delta_{q_1}(E) \leq \frac{15}{16}$ . As  $E$  was arbitrary, we have thus shown  $d(\delta_{q_1}, \delta_{q_2}) \leq \frac{15}{16} < 1$ . We will show in Proposition 12 that we have in fact  $d(\delta_{q_1}, \delta_{q_2}) = \sqrt{2}/4$ , so distances may be irrational. The proof of Proposition 12 shows that  $d(\delta_{q_1}, \delta_{q_2}) = \delta_{q_1}(E) - \delta_{q_2}(E)$  holds for the event

$$E := \{wccc \dots \mid w \in \{a, b\}^*, \#_a(w) \geq \#_b(w)\},$$

where  $\#_a(w)$  and  $\#_b(w)$  denote the number of occurrences of  $a$  and  $b$  in  $w$  respectively. This may be intuitive as  $q_1$  is more likely to emit  $a$ -letters than  $b$ -letters, whereas for  $q_2$  it is the opposite. We remark that this event  $E$  is not  $\omega$ -regular, i.e., it cannot be recognized by a Büchi automaton. As a matter of fact, any  $\omega$ -regular event can only differentiate the two LMCs by a rational number, as the probability of any  $\omega$ -regular event must be rational.

#### 3.2 Example 2

Consider the LMC in Figure 2. Both states  $q_1, q_2$  can initiate *any* run  $r \in \Sigma^\omega$ . Note also that we have  $\delta_{q_1}(\{r\}) = \delta_{q_2}(\{r\}) = 0$  for any single run  $r \in \Sigma^\omega$ . Nevertheless it follows from Theorem 7 that we have  $d(\delta_{q_1}, \delta_{q_2}) = 1$ . Moreover, Theorem 11 will provide an event  $E$  with  $\delta_{q_1}(E) = 1$  and  $\delta_{q_2}(E) = 0$ . Intuitively, such an event could be based on the observation that if  $q_1$  is the initial state, it is more likely after an even number of emitted  $b$ -letters to emit another  $b$ , whereas if  $q_2$  is the initial state, it is more likely after an even number of emitted  $b$ -letters to emit an  $a$ -letter. By the law of large numbers, this difference almost surely “shows” in the long run.

In the following we sketch a proof for the fact that no  $\omega$ -regular event  $E$  satisfies both  $\delta_{q_1}(E) = 1$  and  $\delta_{q_2}(E) = 0$ . In fact, we even show that for any  $\omega$ -regular  $E$  with  $\delta_{q_1}(E) = 1$  we also have  $\delta_{q_2}(E) = 1$ . (We omit precise automata-theoretic definitions here, as this argument will play no further role in this paper.) Let  $E$  be any  $\omega$ -regular event. Let  $R$  be a deterministic Rabin automaton for  $E$ , with initial state  $r_0$ . Let  $\mathcal{M}_R$  denote the LMC obtained by taking the cross-product of  $R$  and the chain from Figure 2. Let  $\delta_{q_1}(E) = 1$ . Then all bottom SCCs of  $\mathcal{M}_R$  reachable from  $(r_0, q_1)$  are accepting. As the qualitative transition structure (i.e., distinguishing only zero and nonzero transition probabilities) is completely symmetric for  $q_1$  and  $q_2$ , it follows that all bottom SCCs of  $\mathcal{M}_R$  reachable from  $(r_0, q_2)$  are accepting as well. Hence we have  $\delta_{q_2}(E) = 1$ .

#### 4. An Approximation Algorithm

In this section we define two computable sequences that converge to the distance from below and from above, respectively. This yields an algorithm for approximating the distance with arbitrary precision.

From now on until the end of Section 5 we fix an LMC  $\mathcal{M} = (Q, \Sigma, M)$  and (initial) distributions  $\pi_1, \pi_2$ . For  $w \in \Sigma^*$  we define

$$\begin{aligned} \min(w) &:= \min\{|\pi_1^w|, |\pi_2^w|\} \quad \text{and} \\ \text{con}(w) &:= \max\{|\mu_1| \mid \mu_1 \leq \pi_1^w \wedge \exists \mu_2 \leq \pi_2^w : \mu_1 \equiv \mu_2\}. \end{aligned}$$

For  $k \in \mathbb{N}$ , we also define  $\min(k) := \sum_{w \in \Sigma^k} \min(w)$  and  $\text{con}(k) := \sum_{w \in \Sigma^k} \text{con}(w)$ . The following proposition lists basic properties of those quantities.

**Proposition 4.** *Let  $w \in \Sigma^*$  and  $k \in \mathbb{N}$ .*

- (a) *We have  $1 \geq \min(w) \geq \text{con}(w) \geq 0$ . Hence  $\min(k) \geq \text{con}(k)$ .*
- (b) *We have  $\min(w) \geq \sum_{a \in \Sigma} \min(wa)$  and  $\text{con}(w) \leq \sum_{a \in \Sigma} \text{con}(wa)$ . Hence we have  $\min(k) \geq \min(k+1)$  and  $\text{con}(k) \leq \text{con}(k+1)$ .*
- (c) *The limits  $\min(\infty) := \lim_{i \rightarrow \infty} \min(i)$  and  $\text{con}(\infty) := \lim_{i \rightarrow \infty} \text{con}(i)$  exist, and we have  $\min(\infty) \geq \text{con}(\infty)$ .*

*Proof.*

- (a) We have  $1 \geq |\pi_1^w| \geq \min\{|\pi_1^w|, |\pi_2^w|\} = \min(w)$ , hence  $1 \geq \min(w)$ . Clearly,  $\text{con}(w) \geq 0$ .  
Let  $\mu_1, \mu_2$  be the subdistributions such that  $\text{con}(w) = |\mu_1|$  and  $\pi_1^w \geq \mu_1 \equiv \mu_2 \leq \pi_2^w$ . Then we have  $|\pi_1^w| \geq |\mu_1| = |\mu_2| \leq |\pi_2^w|$ , hence  $\text{con}(w) = |\mu_1| \leq \min\{|\pi_1^w|, |\pi_2^w|\} = \min(w)$ .
- (b) Let  $i \in \{1, 2\}$  with  $|\pi_i^w| \leq |\pi_{3-i}^w|$ . Then we have:

$$\begin{aligned} \min(w) &= \min\{|\pi_1^w|, |\pi_2^w|\} = |\pi_i^w| = |\pi_i M(w)| \\ &= |\pi_i M(w) \sum_{a \in \Sigma} M(a)| \\ &= \sum_{a \in \Sigma} |\pi_i M(wa)| = \sum_{a \in \Sigma} |\pi_i^{wa}| \\ &\geq \sum_{a \in \Sigma} \min\{|\pi_i^{wa}|, |\pi_{3-i}^{wa}|\} = \sum_{a \in \Sigma} \min(wa) \end{aligned}$$

Let  $\mu_1, \mu_2$  be the subdistributions such that  $\text{con}(w) = |\mu_1|$  and  $\pi_1^w \geq \mu_1 \equiv \mu_2 \leq \pi_2^w$ . It follows that, for all  $a \in \Sigma$ , we have  $\pi_1^{wa} \geq \mu_1 M(a) \equiv \mu_2 M(a) \leq \pi_2^{wa}$ , hence  $\text{con}(wa) \geq |\mu_1 M(a)|$ . So we have:

$$\begin{aligned} \sum_{a \in \Sigma} \text{con}(wa) &\geq \sum_{a \in \Sigma} |\mu_1 M(a)| = |\mu_1| \sum_{a \in \Sigma} M(a) = |\mu_1| \\ &= \text{con}(w). \end{aligned}$$

- (c) Follows from (a) and (b).  $\square$

The quantities  $\min(k)$  and  $\text{con}(k)$  provide lower and upper bounds for the distance:

**Proposition 5.** *For all  $k \in \mathbb{N}$  we have:*

$$1 - \min(k) \leq d(\pi_1, \pi_2) \leq 1 - \text{con}(k).$$

*Proof.* We show first the lower bound. Let  $k \in \mathbb{N}$ . Define  $W_1 := \{w \in \Sigma^k \mid |\pi_1^w| \geq |\pi_2^w|\}$  and  $W_2 := \{w \in \Sigma^k \mid |\pi_1^w| < |\pi_2^w|\}$ .

By the definitions we have:

$$\begin{aligned} d(\pi_1, \pi_2) &= \sup_{E \subseteq \Sigma^\omega} (\pi_1(E) - \pi_2(E)) \\ &\geq \pi_1(W_1 \Sigma^\omega) - \pi_2(W_1 \Sigma^\omega) \\ &= 1 - \pi_1(W_2 \Sigma^\omega) - \pi_2(W_1 \Sigma^\omega) \\ &= 1 - \sum_{w \in W_2} |\pi_1^w| - \sum_{w \in W_1} |\pi_2^w| \\ &= 1 - \sum_{w \in \Sigma^k} \min\{|\pi_1^w|, |\pi_2^w|\} \\ &= 1 - \min(k). \end{aligned}$$

Now we show the upper bound. For an event  $E \subseteq \Sigma^\omega$  and a word  $w \in \Sigma^*$ , we denote by  $w^{-1}E$  the event  $\{u \in \Sigma^\omega \mid wu \in E\}$ . For  $w \in \Sigma^*$  we write  $\mu_1^{(w)}$  and  $\mu_2^{(w)}$  to denote subdistributions with  $\text{con}(w) = |\mu_1^{(w)}| = |\mu_2^{(w)}|$  and  $\pi_1^w \geq \mu_1^{(w)} \equiv \mu_2^{(w)} \leq \pi_2^w$ . The following inequalities hold:

$$\begin{aligned} \pi_1^w(w^{-1}E) &= \mu_1^{(w)}(w^{-1}E) + (\pi_1^w - \mu_1^{(w)})(w^{-1}E) \\ &\leq \mu_1^{(w)}(w^{-1}E) + |\pi_1^w| - |\mu_1^{(w)}| \\ \pi_2^w(w^{-1}E) &\geq \mu_2^{(w)}(w^{-1}E) \end{aligned} \quad (3)$$

We have:

$$\begin{aligned} d(\pi_1, \pi_2) &= \sup_{E \subseteq \Sigma^\omega} (\pi_1(E) - \pi_2(E)) \\ &= \sup_{E \subseteq \Sigma^\omega} \sum_{w \in \Sigma^k} \pi_1^w(w^{-1}E) - \pi_2^w(w^{-1}E) \\ &\leq \sup_{E \subseteq \Sigma^\omega} \sum_{w \in \Sigma^k} \mu_1^{(w)}(w^{-1}E) + |\pi_1^w| - |\mu_1^{(w)}| \\ &\quad - \mu_2^{(w)}(w^{-1}E) \quad (\text{by (3)}) \\ &= \sup_{E \subseteq \Sigma^\omega} \sum_{w \in \Sigma^k} |\pi_1^w| - |\mu_1^{(w)}| \quad (\text{as } \mu_1^{(w)} \equiv \mu_2^{(w)}) \\ &= 1 - \sum_{w \in \Sigma^k} |\mu_1^{(w)}| = 1 - \text{con}(k) \end{aligned}$$

$\square$

The lower bound in this proposition follows by considering the event  $E_k := W_1 \Sigma^\omega$  (where  $W_1$  is from the proof), which depends only on the length- $k$  prefix of the run. In fact, if we restrict each run to its length- $k$  prefix, we obtain a finite sample space, and the event  $E_k$  is the maximizer according to (2) in the introduction. We could define, for each  $k \in \mathbb{N}$ , a distance  $d_k(\pi_1, \pi_2)$  with

$$\begin{aligned} d_k(\pi_1, \pi_2) &= \max_{W \in \Sigma^k} |\pi_1(W \Sigma^\omega) - \pi_2(W \Sigma^\omega)| \\ &= \pi_1(E_k) - \pi_2(E_k) = 1 - \min(k). \end{aligned}$$

Since  $d_k(\pi_1, \pi_2) \leq d_{k+1}(\pi_1, \pi_2)$  holds by Proposition 4 (b), there is a limit  $\lim_{k \rightarrow \infty} d_k(\pi_1, \pi_2)$ , which equals  $d(\pi_1, \pi_2)$  (as we will show in Theorem 7). This would offer an alternative but equivalent definition of the distance, which avoids the use of infinite runs by replacing them with increasing prefixes.

By combining Propositions 4 and 5 we obtain

$$1 - \min(\infty) \leq d(\pi_1, \pi_2) \leq 1 - \text{con}(\infty). \quad (4)$$

In the rest of this section we show that those inequalities are in fact equalities.

Recall that for a (random) run  $r \in \Sigma^\omega$  we write  $r_i \in \Sigma^i$  for the length- $i$  prefix of  $r$ . For  $i \in \mathbb{N}$ , we define the random variable  $L_i$

that assigns to a run  $r \in \Sigma^\omega$  the likelihood ratio

$$L_i(r) := |\pi_2^{r_i}|/|\pi_1^{r_i}|.$$

Observe that  $L_0(r) = 1$ .

**Proposition 6.** *We have*

$$\begin{aligned} \pi_1 \left( \lim_{i \rightarrow \infty} L_i \text{ exists and is in } [0, \infty) \right) &= 1 \quad \text{and} \\ \pi_2 \left( \lim_{i \rightarrow \infty} 1/L_i \text{ exists and is in } [0, \infty) \right) &= 1. \end{aligned}$$

*Proof.* We prove only the first equality; the second equality is proved similarly. First we show that the sequence  $L_0, L_1, \dots$  is a martingale. Denote by  $\text{Ex}_1$  the expectation with respect to  $\pi_1$ . Let  $i \in \mathbb{N}$  and let  $w \in \Sigma^i$  with  $|\pi_1^w| > 0$ . We have:

$$\begin{aligned} \text{Ex}_1(L_{i+1} \mid \text{Run}(w)) &= \sum_{q, q' \in Q} \sum_{a \in \Sigma} \frac{\pi_1^w(q) M(a)(q, q')}{|\pi_1^w|} \cdot \frac{|\pi_2^{wa}|}{|\pi_1^{wa}|} \\ &= \frac{1}{|\pi_1^w|} \sum_{a \in \Sigma} \frac{|\pi_2^{wa}|}{|\pi_1^{wa}|} \underbrace{\sum_{q \in Q} \pi_1^w(q) \sum_{q' \in Q} M(a)(q, q')}_{= |\pi_1^{wa}|} \\ &= \frac{1}{|\pi_1^w|} \sum_{a \in \Sigma} |\pi_2^{wa}| = |\pi_2^w|/|\pi_1^w| = L_i \end{aligned}$$

So  $L_0, L_1, \dots$  is a martingale. More precisely, the sequence  $L_0, L_1, \dots$  is a nonnegative martingale with  $\text{Ex}_1(L_i) = 1$  for all  $i \in \mathbb{N}$ . So the martingale convergence theorem (more precisely, “Doob’s forward convergence theorem”, see e.g. [20]) applies, and we obtain  $\pi_1(\lim_{i \rightarrow \infty} L_i \text{ exists and is finite}) = 1$ .  $\square$

In the following we may write  $\lim_{i \rightarrow \infty} L_i = \infty$  to mean  $\lim_{i \rightarrow \infty} 1/L_i = 0$ . Define

$$\bar{L} := \lim_{i \rightarrow \infty} L_i \in [0, \infty] \quad (\text{if the limit exists}). \quad (5)$$

The random variable  $\bar{L}$  plays a crucial role in the next section and is also used in the proof of the following theorem.

**Theorem 7.** *We have*

$$1 - \min(\infty) = d(\pi_1, \pi_2) = 1 - \text{con}(\infty).$$

*Proof sketch.* The proof (Appendix B) is somewhat technical and we only give a sketch here. Considering (4) it suffices to show that  $\min(\infty) = \text{con}(\infty)$ . By Proposition 4 we have  $\min(k) \geq \text{con}(k)$  for all  $k$ , so loosely speaking we have to show that for “large”<sup>1</sup>  $k$ ,  $\min(k)$  is not much larger than  $\text{con}(k)$ . We first show that this holds for individual runs started from  $\pi_1$ ; more precisely, we show for all  $\gamma > 0$  that

$$\begin{aligned} \pi_1(\bar{L} > 0) \\ = \pi_1(\bar{L} > 0 \wedge \exists i \in \mathbb{N} : \min(r_i) \leq (1 + \gamma) \text{con}(r_i)). \end{aligned} \quad (6)$$

In words: Conditioned under the event  $\{\bar{L} > 0\}$  the probability that eventually  $\min(r_i) \leq (1 + \gamma) \text{con}(r_i)$  holds is 1. To show (6) we first show that conditioned under  $\{\bar{L} > 0\}$  we have with probability 1 that the distance between the distributions  $\pi_1^{r_i}/|\pi_1^{r_i}|$  and  $\pi_2^{r_i}/|\pi_2^{r_i}|$  converges to 0. Using the fact that the set of distributions is compact, one can then show (6).

To show that for large  $k$ ,  $\min(k)$  is not much larger than  $\text{con}(k)$ , we consider a partition  $\Sigma^k = W_1 \cup W_2 \cup W_3$ . The set  $W_1$  contains the words  $w$  with small  $|\pi_2^w|/|\pi_1^w|$ . So

$\sum_{w \in W_1} \min(w) \leq \sum_{w \in W_1} |\pi_2^w|$  is small. The set  $W_2$  contains the words  $w$  with  $\min(w) > (1 + \gamma) \text{con}(w)$  and large  $|\pi_2^w|/|\pi_1^w|$ . Runs with prefixes in  $W_2$  and  $\bar{L} = 0$  are unlikely, as  $L_k = |\pi_2^w|/|\pi_1^w|$  is large and  $k$  is large and the sequence  $L_0, L_1, \dots$  converges to  $\bar{L}$  by Proposition 6. Runs with prefixes in  $W_2$  and  $\bar{L} > 0$  are also unlikely because of (6). So  $\sum_{w \in W_2} \min(w) \leq \sum_{w \in W_2} |\pi_1^w|$  is small. Finally, the set  $W_3$  contains the words  $w$  with  $\min(w) \leq (1 + \gamma) \text{con}(w)$  and large  $|\pi_2^w|/|\pi_1^w|$ . So  $\sum_{w \in W_3} \min(w)$  is (for small  $\gamma$ ) not much larger than  $\sum_{w \in W_3} \text{con}(w) \leq \text{con}(k)$ . By adding the mentioned inequalities we obtain that  $\min(k) = \sum_{w \in \Sigma^k} \min(w)$  is not much larger than  $\text{con}(k)$ .  $\square$

**Corollary 8.** *There is an algorithm that, given  $\varepsilon > 0$ , computes  $a \in \mathbb{Q}$  such that  $d(\pi_1, \pi_2) \in [a, a + \varepsilon]$ .*

*Proof.* By Proposition 5 and Theorem 7 the sequences  $(1 - \min(k))_{k \in \mathbb{N}}$  and  $(1 - \text{con}(k))_{k \in \mathbb{N}}$  converge to  $d(\pi_1, \pi_2)$  from below and above, respectively. For each  $k$ , the values  $\min(k)$  and  $\text{con}(k)$  are computable.  $\square$

In terms of the complexity of approximating the distance we have the following result:

**Proposition 9.** *Approximating the distance up to any  $\varepsilon$  whose size is polynomial in the given LMC is NP-hard with respect to Turing reductions.*

*Proof.* In [14, Section 6] (see also [4, Theorem 7]), a reduction is given from the *clique decision problem* to show that computing the distance in LMCs is NP-hard. In their reduction the distance is rational and of polynomial size in the input. Using the continued-fraction method (see e.g. Section 2.4 of [7] for an explanation) it follows that a polynomial-time algorithm (if it exists) for approximating the distance can be used to construct a polynomial-time algorithm for computing the distance exactly. Hence the conclusion.  $\square$

This NP-hardness result also follows from the proof of [4, Theorem 10].

## 5. A Maximizing Event

The proof of Theorem 7 does not yield an event  $E_1$  with  $\pi_1(E_1) - \pi_2(E_1) = d(\pi_1, \pi_2) \stackrel{\text{def}}{=} \sup_{E \subseteq \Sigma^\omega} |\pi_1(E) - \pi_2(E)|$ . In fact, it is not clear a priori whether such an event exists. In this section we exhibit such a “witness”  $E_1$ . It follows that the supremum from the definition of distance is in fact a maximum.

For some intuition recall from (2) in the introduction that in the countable case the event  $E_1 = \{r \in \Omega \mid \pi_1(r) \geq \pi_2(r)\}$  is the desired maximizer. In the case of LMCs this does not work, since each individual run may have probability 0 (as, e.g., in Figure 2). However, by rewriting the inequality  $\pi_1(r) \geq \pi_2(r)$  as  $\pi_2(r)/\pi_1(r) \leq 1$ , one is tempted to guess that  $\pi_2(r)/\pi_1(r)$  can be replaced by  $\bar{L}(r)$  as defined in (5). In the rest of the section we show that this intuition is correct. Define the events

$$E_1 := \{\bar{L} \leq 1\} \quad \text{and} \quad E_2 := \{\bar{L} > 1\}.$$

By Proposition 6 we have

$$\pi_1(E_1) + \pi_1(E_2) = \pi_2(E_1) + \pi_2(E_2) = 1. \quad (7)$$

The following lemma will suffice for showing that  $E_1$  is the desired maximizer.

**Lemma 10.** *We have  $\pi_1(E_2) + \pi_2(E_1) \leq \min(\infty)$ .*

<sup>1</sup>In the rest of this proof sketch we gloss over the precise meaning of “small”, “not much larger”, etc., and omit the quotation marks.

*Proof.* Towards a contradiction, suppose that this does not hold. Then there is  $k' \in \mathbb{N}$  with  $\pi_1(E_2) + \pi_2(E_1) > \min(k')$ ; hence there is  $\gamma > 0$  with

$$\pi_1(E_2) + \pi_2(E_1) > \min(k') + 4\gamma. \quad (8)$$

Choose  $\varepsilon \in (0, \gamma]$  small enough so that

$$\pi_1(\bar{L} \in (1, 1 + 2\varepsilon]) \leq \gamma.$$

Using Proposition 6, choose  $k \geq k'$  large enough so that we have

$$\begin{aligned} \pi_1(L_k \leq 1 + \varepsilon \wedge \bar{L} > 1 + 2\varepsilon) &\leq \gamma \quad \text{and} \\ \pi_2(L_k > 1 + \varepsilon \wedge \bar{L} \leq 1) &\leq \gamma. \end{aligned} \quad (9)$$

Then we have:

$$\begin{aligned} \pi_1(E_2) &= \pi_1(\bar{L} \in (1, 1 + 2\varepsilon]) + \pi_1(\bar{L} > 1 + 2\varepsilon) \quad (\text{def. of } E_2) \\ &\leq \gamma + \pi_1(\bar{L} > 1 + 2\varepsilon) \quad (\text{choice of } \varepsilon) \\ &= \gamma + \pi_1(L_k > 1 + \varepsilon \wedge \bar{L} > 1 + 2\varepsilon) \\ &\quad + \pi_1(L_k \leq 1 + \varepsilon \wedge \bar{L} > 1 + 2\varepsilon) \\ &\leq \gamma + \pi_1(L_k > 1 + \varepsilon) \\ &\quad + \pi_1(L_k \leq 1 + \varepsilon \wedge \bar{L} > 1 + 2\varepsilon) \\ &\leq 2\gamma + \pi_1(L_k > 1 + \varepsilon) \quad (\text{by (9)}) \end{aligned}$$

Similarly we have:

$$\begin{aligned} \pi_2(E_1) &\leq \pi_2(L_k > 1 + \varepsilon \wedge \bar{L} \leq 1) \\ &\quad + \pi_2(L_k \leq 1 + \varepsilon) \quad (\text{def. of } E_1) \\ &\leq \gamma + \pi_2(L_k \leq 1 + \varepsilon) \quad (\text{by (9)}) \end{aligned}$$

By adding those two inequalities we obtain

$$\begin{aligned} \pi_1(E_2) + \pi_2(E_1) &\leq 3\gamma + \pi_1(L_k > 1 + \varepsilon) + \pi_2(L_k \leq 1 + \varepsilon). \end{aligned} \quad (10)$$

Define the partition of  $\Sigma^k$  in  $\Sigma^k = W_1 \cup W_2 \cup W_3$  with

$$\begin{aligned} W_1 &:= \{w \in \Sigma^k \mid |\pi_2^w|/|\pi_1^w| \leq 1\} \\ W_2 &:= \{w \in \Sigma^k \mid 1 < |\pi_2^w|/|\pi_1^w| \leq 1 + \varepsilon\} \\ W_3 &:= \{w \in \Sigma^k \mid 1 + \varepsilon < |\pi_2^w|/|\pi_1^w|\}. \end{aligned}$$

Then we have

$$\begin{aligned} \pi_2(L_k \leq 1) &= \sum_{w \in W_1} |\pi_2^w| \\ \pi_2(1 < L_k \leq 1 + \varepsilon) &= \sum_{w \in W_2} |\pi_2^w| \leq \sum_{w \in W_2} (1 + \varepsilon) |\pi_1^w| \\ \pi_1(L_k > 1 + \varepsilon) &= \sum_{w \in W_3} |\pi_1^w|. \end{aligned}$$

By adding those (in)equalities we obtain

$$\begin{aligned} \pi_1(L_k > 1 + \varepsilon) + \pi_2(L_k \leq 1 + \varepsilon) &\leq (1 + \varepsilon) \min(k) \leq \min(k) + \varepsilon \leq \min(k) + \gamma. \end{aligned}$$

Combining this with (10) yields

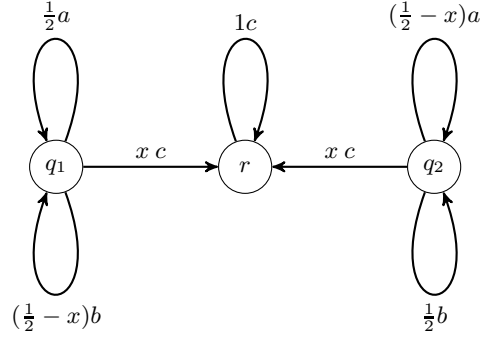
$$\pi_1(E_2) + \pi_2(E_1) \leq \min(k) + 4\gamma \leq \min(k') + 4\gamma,$$

thus contradicting (8) as desired.  $\square$

Now we can prove that  $E_1$  is the desired maximizing event.

**Theorem 11.** *We have*

$$d(\pi_1, \pi_2) = \pi_1(E_1) - \pi_2(E_1).$$



**Figure 3.** In this LMC,  $x \in (0, \frac{1}{2})$  is a parameter. For  $x = \frac{1}{4}$  the Dirac distributions  $\delta_{q_1}$  and  $\delta_{q_2}$  have distance  $\sqrt{2}/4 \notin \mathbb{Q}$ .

*Proof.* We have:

$$\begin{aligned} d(\pi_1, \pi_2) &\geq \pi_1(E_1) - \pi_2(E_1) \quad (\text{definition of distance}) \\ &= 1 - \pi_1(E_2) - \pi_2(E_1) \quad (\text{by (7)}) \\ &\geq 1 - \min(\infty) \quad (\text{Lemma 10}) \\ &= d(\pi_1, \pi_2) \quad (\text{Theorem 7}) \end{aligned}$$

$\square$

## 6. Irrational Distances and Lower Bounds

The following proposition shows that the distance can be irrational even if all numbers in the description of the LMC are rational.

**Proposition 12.** *Consider the LMC shown in Figure 3, with parameter  $x \in (0, \frac{1}{2})$ . We have  $d(\delta_{q_1}, \delta_{q_2}) = \frac{1}{2}\sqrt{2x}$ .*

We start with a technical lemma.

**Lemma 13.** *For  $y \in [0, \frac{1}{4})$  we have*

$$\sum_{n=0}^{\infty} \binom{2n}{n} y^n = \frac{1}{\sqrt{1-4y}}.$$

*Proof.* By a binomial series we have:

$$\frac{1}{\sqrt{1-4y}} = (1-4y)^{-1/2} = \sum_{n=0}^{\infty} \binom{-1/2}{n} (-4y)^n$$

By induction on  $n \in \mathbb{N}$  one can show that  $\binom{-1/2}{n} = \frac{(2n)}{n} \frac{(-1)^n}{4^n}$ . The lemma follows.  $\square$

*Proof of Proposition 12.* We write  $\pi_1 := \delta_{q_1}$  and  $\pi_2 := \delta_{q_2}$ . Define  $C := \{wccc \dots \mid w \in \{a, b\}^*\} \subseteq \Sigma^\omega$ . Clearly we have  $\pi_1(C) = \pi_2(C) = 1$ . Define

$$\begin{aligned} E_{>} &:= \{wccc \dots \mid w \in \Sigma^*, \#_a(w) > \#_b(w)\} \quad \text{and} \\ E_{=} &:= \{wccc \dots \mid w \in \Sigma^*, \#_a(w) = \#_b(w)\}, \end{aligned}$$

where  $\#_a(w)$  and  $\#_b(w)$  denote the number of occurrences of  $a$  resp.  $b$  in the word  $w$ . The events  $E_{>}, E_{<}, E_{\leq}$  are defined accordingly.

Recall the event  $E_1 = \{\bar{L} \leq 1\} \subseteq \Sigma^\omega$  from Section 5. Using the fact that the LMC in Figure 3 is “deterministic” (i.e., for each  $a \in \Sigma$  and  $q \in Q$  there is at most one  $q' \in Q$  with  $M(a)(q, q') > 0$ ), it is easy to verify that we have  $E_1 \cap C = E_{>}$ .

We have:

$$\begin{aligned}
d(\pi_1, \pi_2) &= \pi_1(E_1) - \pi_2(E_1) && \text{(Theorem 11)} \\
&= \pi_1(E_1 \cap C) - \pi_2(E_1 \cap C) && (\text{as } \pi_1(C) = \pi_2(C) = 1) \\
&= \pi_1(E_{\geq}) - \pi_2(E_{\geq}) && (\text{as argued above}) \\
&= \pi_1(E_{\geq}) - \pi_1(E_{\leq}) && (\text{by symmetry of the chain}) \\
&= \pi_1(E_{\geq}) - (1 - \pi_1(E_{>})) && (\text{as } \pi_1(C) = 1) \\
&= 2\pi_1(E_{>}) + \pi_1(E_{=}) - 1 && (\text{by the definitions}). \quad (11)
\end{aligned}$$

The following identity is proved in [10, p.167, (5.20)] and in [13] with a short combinatorial proof:

$$\sum_{m=0}^n \binom{m+n}{m} \left(\frac{1}{2}\right)^m = 2^n \quad \text{for } n \in \mathbb{N}. \quad (12)$$

For  $m, n \in \mathbb{N}$  define  $E(m, n) := \{wccc \dots \mid w \in \Sigma^*, \#_a(w) = m, \#_b(w) = n\}$ . We have:

$$\begin{aligned}
\pi_1(E_{\leq}) &= \sum_{n=0}^{\infty} \sum_{m=0}^n \pi_1(E(m, n)) \\
&= \sum_{n=0}^{\infty} \underbrace{\sum_{m=0}^n \binom{m+n}{m} \left(\frac{1}{2}\right)^m \left(\frac{1}{2} - x\right)^n}_{=2^n \text{ by (12)}} x \\
&= \frac{x}{1 - 2(\frac{1}{2} - x)} = \frac{1}{2}
\end{aligned}$$

So we have  $\pi_1(E_{>}) = 1 - \pi_1(E_{\leq}) = \frac{1}{2}$  and hence by (11)

$$d(\pi_1, \pi_2) = \pi_1(E_{=}). \quad (13)$$

We have:

$$\begin{aligned}
\pi_1(E_{=}) &= \sum_{n=0}^{\infty} \pi_1(E(n, n)) \\
&= \sum_{n=0}^{\infty} \binom{2n}{n} \left(\frac{1}{2}\right)^n \left(\frac{1}{2} - x\right)^n x \\
&= \frac{x}{\sqrt{1 - 4(\frac{1}{4} - \frac{1}{2}x)}} = \frac{1}{2} \sqrt{2x} \quad (\text{by Lemma 13}).
\end{aligned}$$

so the statement follows with (13).  $\square$

Note that when  $x = \frac{1}{4}$ , the LMC shown in Figure 3 is essentially the union of the two LMCs shown in Figure 1. Proposition 12 states that  $d(\delta_{q_1}, \delta_{q_2}) = \sqrt{2}/4$ , thus substantiating a claim in Section 3.1.

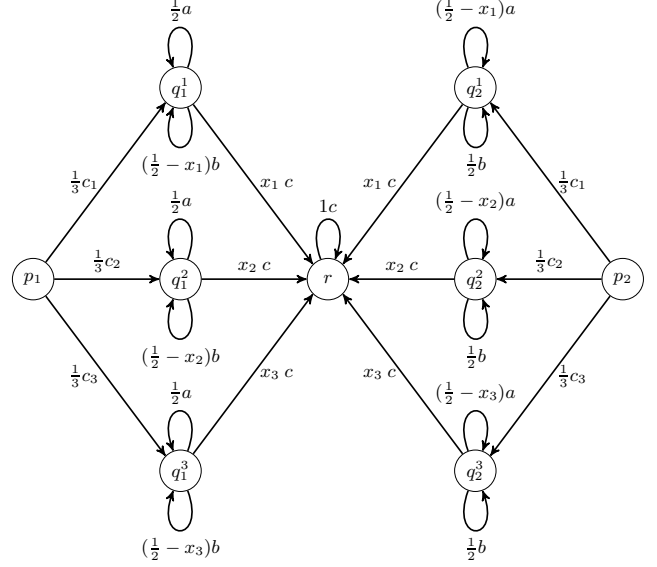
This example suggests that in general it is not obvious what *computing* the distance means, as it may be irrational. Nevertheless it is shown in [14, Section 6] that computing the distance is NP-hard (with respect to Turing reductions). In that reduction the computed LMCs have a rational distance by construction. However, in light of Proposition 12 it may be more natural to study the *threshold-distance* problem defined as follows: Given an LMC, two initial distributions  $\pi_1, \pi_2$ , and a threshold  $\tau \in [0, 1] \cap \mathbb{Q}$ , decide whether  $d(\pi_1, \pi_2) \geq \tau$ .

By Proposition 9, together with a binary search, the following lower bound follows:

**Proposition 14.** *The threshold-distance problem is NP-hard with respect to Turing reductions.*

We remark that this can also be done by modifying the reduction from [14], see Appendix C.

In the following we give another lower bound for the threshold-distance problem: the problem is hard for the square-root-sum



**Figure 4.** This LMC is obtained by combining the chain from Figure 3 in parallel  $n = 3$  times. We have  $d(\delta_{p_1}, \delta_{p_2}) = \frac{1}{3} (d(\delta_{q_1^1}, \delta_{q_2^1}) + d(\delta_{q_1^2}, \delta_{q_2^2}) + d(\delta_{q_1^3}, \delta_{q_2^3}))$ .

problem, as we explain now. Following [1] the *square-root-sum* problem is defined as follows. Given natural numbers  $s_1, \dots, s_n \in \mathbb{N}$  and  $t \in \mathbb{N}$ , decide whether  $\sum_{i=1}^n \sqrt{s_i} \geq t$ . Membership of square-root-sum in NP has been open since 1976 when Garey, Graham and Johnson [8] showed NP-hardness of the travelling-salesman problem with Euclidean distances, but left membership in NP open. It is known that square-root-sum reduces to PosSLP and hence lies in the 4th level of the counting hierarchy, see [1] and the references therein for more information on square-root-sum, PosSLP, and the counting hierarchy.

We use the LMC from Figure 3 as a “gadget” to prove hardness for the square-root-sum problem:

**Theorem 15.** *There is a polynomial-time many-one reduction from the square-root-sum problem to the threshold-distance problem.*

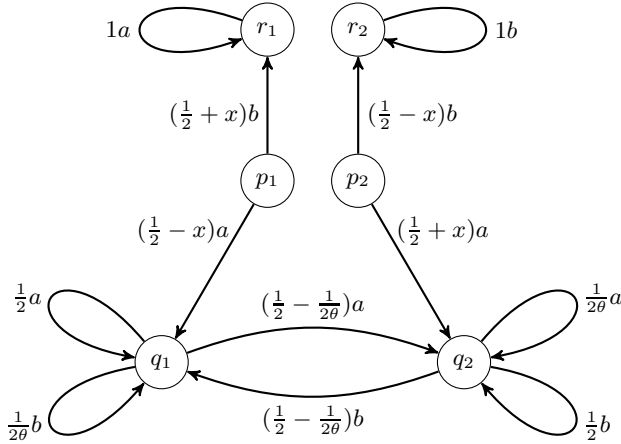
*Proof sketch.* The construction is by taking the LMC from Figure 3 as a gadget, and joining  $n$  instances of it in parallel. This is sketched for  $n = 3$  in Figure 4. In general we have  $\Sigma = \{c_1, \dots, c_n, a, b, c\}$  and  $Q = \{p_1, p_2, q_1^1, \dots, q_1^n, q_2^1, \dots, q_2^n, r\}$ . Using this construction we have

$$d(\delta_{p_1}, \delta_{p_2}) = \frac{1}{n} \sum_{i=1}^n d(\delta_{q_1^i}, \delta_{q_2^i}). \quad (14)$$

We prove (14) in Appendix C. From the proof of Proposition 12 we know the distances  $d(\delta_{q_1^i}, \delta_{q_2^i})$  and the corresponding maximizing events. The proof is completed by suitably choosing the  $x_i$  and the threshold  $\tau$ , see Appendix C.  $\square$

### 6.1 Bernoulli Convolutions

In this section, we establish another “lower bound” by demonstrating a link to Bernoulli convolutions. Consider the LMC in Figure 5 which has two parameters:  $\theta > 1$  and  $x \in [-\frac{1}{2}, \frac{1}{2}]$ . For each  $\theta > 1$ , denote by  $d_\theta : [-\frac{1}{2}, \frac{1}{2}] \rightarrow [0, 1]$  the function such that  $d_\theta(x)$  is the distance between states  $p_1$  and  $p_2$  in the chain with parameters  $\theta$  and  $x$ . Using the Banach fixed-point theorem one can show (see Appendix D):



**Figure 5.** The distance between state  $p_1, p_2$  depends on a Bernoulli-convolution.

**Proposition 16.** For all  $\theta > 1$  we have  $d_\theta(x) = \frac{1}{2} + \frac{1}{2}f_\theta(x)$  for the unique function  $f_\theta : \mathbb{R} \rightarrow \mathbb{R}$  with

$$f_\theta(x) = \begin{cases} -2x & x \leq -\frac{1}{2} \\ \frac{1}{2\theta}f_\theta(\theta x - (\frac{1}{2}\theta - \frac{1}{2})) & x \in [-\frac{1}{2}, \frac{1}{2}] \\ +\frac{1}{2\theta}f_\theta(\theta x + (\frac{1}{2}\theta - \frac{1}{2})) & x \geq +\frac{1}{2} \end{cases}$$

It follows that the derivative of  $f_\theta$  must satisfy

$$f'_\theta(x) = \begin{cases} -2 & x \leq -\frac{1}{2} \\ \frac{1}{2\theta}f'_\theta(\theta x - (\frac{1}{2}\theta - \frac{1}{2})) & x \in [-\frac{1}{2}, \frac{1}{2}] \\ +\frac{1}{2\theta}f'_\theta(\theta x + (\frac{1}{2}\theta - \frac{1}{2})) & x \geq +\frac{1}{2} \end{cases} \quad (15)$$

Again, one can use the Banach fixed-point theorem to show that the solution  $f'_\theta$  is unique.

The functional equation (15) is known from the study of *Bernoulli convolutions*, see [16] for a survey and [2, Chapter 5] for a gentle introduction. In this field the solution of (15) occurs (translated and rescaled) as the cumulative distribution function of the random variable  $\sum_{i=0}^{\infty} X_i/\theta^i$ , where the  $X_i$  are random variables that take on  $-1$  and  $+1$  with probability  $\frac{1}{2}$  each. Bernoulli convolutions have been studied since the 1930s. It is known that the solutions of (15) are either absolutely continuous or singular on  $[-\frac{1}{2}, \frac{1}{2}]$ , depending on  $\theta$ . For  $\theta > 2$  they are singular; in fact, for  $\theta = 3$  the function is the (ternary) Cantor function. For  $\theta = 2$  we have  $f'_\theta(x) = 4x$  for  $x \in [-\frac{1}{2}, \frac{1}{2}]$ . Erdős showed that if  $\theta$  is a *Pisot number*<sup>2</sup>, then  $f'_\theta$  is singular. However, for almost all  $\theta \in (1, 2]$  the function  $f'_\theta$  is absolutely continuous. It is open, e.g., for  $\theta = 3/2$  whether  $f'_\theta$  is absolutely continuous or purely singular.

We conclude from this relation to Bernoulli convolutions that the distance can depend on the probabilities in the LMC in intricate ways.

<sup>2</sup> A Pisot number is a real algebraic integer greater than 1 such that all its Galois conjugates are less than 1 in absolute value. The smallest Pisot number ( $\approx 1.3247$ ) is the real root of  $x^3 - x - 1$ . Another one is the golden ratio  $(\sqrt{5} + 1)/2 \approx 1.6180$ .

## 7. The Distance-1 Problem

The *distance-1 problem* asks whether  $d(\pi_1, \pi_2) = 1$  holds for a given LMC and two distributions  $\pi_1, \pi_2$ . For the rest of the section we fix an LMC  $\mathcal{M} = (Q, \Sigma, M)$  and initial distributions  $\pi_1, \pi_2$ . Recall from Proposition 3 that  $d(\pi_1, \pi_2) = 0$  is equivalent to  $\pi_1 \equiv \pi_2$ , and that the latter problem, language equivalence, is known to be decidable in polynomial time [12]. In this section we show that the distance-1 problem can also be decided in polynomial time. The algorithm and its correctness argument are much more subtle. The following proposition provides a characterisation of the case  $d(\pi_1, \pi_2) < 1$ .

**Proposition 17.** We have  $d(\pi_1, \pi_2) < 1$  if and only if there are  $w \in \Sigma^*$  and subdistributions  $\mu_1, \mu_2$  with  $\mu_1 \leq \pi_1^w$  and  $\mu_2 \leq \pi_2^w$  and  $\mu_1 \equiv \mu_2$  and  $|\mu_1| = |\mu_2| > 0$ .

Note that  $\mu_1 \equiv \mu_2$  implies  $|\mu_1| = |\mu_2|$ . Proposition 17 follows immediately from Theorem 7.

Given  $\pi_1, \pi_2$  and a word  $w \in \Sigma^*$  one can compute  $\pi_1^w$  and  $\pi_2^w$  in polynomial time. Consider the following condition on  $w$ :

$$\exists \mu_1, \mu_2 : \mu_1 \leq \pi_1^w \text{ and } \mu_2 \leq \pi_2^w \text{ and } \mu_1 \equiv \mu_2 \text{ and } |\mu_1| > 0. \quad (16)$$

By Proposition 3 (b), (16) amounts to a feasibility test of a linear program, and hence can be decided in polynomial time. By Proposition 17 we have  $d(\pi_1, \pi_2) < 1$  if and only if there is  $w \in \Sigma^*$  such that (16) holds.

For notational convenience we write  $\text{supp}(w)$  for the pair  $(\text{supp}(\pi_1^w), \text{supp}(\pi_2^w))$  in the following. The condition (16) on  $w$  is in fact only a condition on  $\text{supp}(w)$ , as  $\mu_1 \equiv \mu_2$  implies  $a\mu_1 \equiv a\mu_2$  for all  $a \in [0, \infty)$ . So (16) can be rephrased as

$$\exists \mu_1, \mu_2 : \text{supp}(\mu_1) \subseteq \text{supp}(\pi_1^w) \text{ and } \text{supp}(\mu_2) \subseteq \text{supp}(\pi_2^w) \text{ and } \mu_1 \equiv \mu_2 \text{ and } |\mu_1| > 0. \quad (17)$$

Moreover, for any two words  $w, w' \in \Sigma^*$  with  $\text{supp}(w) = \text{supp}(w')$  we have  $\text{supp}(wa) = \text{supp}(w'a)$  for all  $a \in \Sigma$ . This implies

$$\{\text{supp}(w) \mid w \in \Sigma^*\} = \{\text{supp}(w) \mid w \in \Sigma^*, |w| \leq 2^{2|Q|}\}.$$

This suggests the following nondeterministic algorithm for checking whether  $d(\pi_1, \pi_2) < 1$  holds: compute  $\text{supp}(w)$  for a guessed word  $w \in \Sigma^*$  with  $|w| \leq 2^{2|Q|}$  and check (17) for feasibility. Note that  $w$  may have exponential length but need not be stored as a whole. This results in a PSPACE algorithm.

In the following, we give a polynomial-time algorithm, which is based on further properties of the distance.

Given subdistributions  $\mu_1, \mu_2$  with  $|\mu_1|, |\mu_2| > 0$  we define the following relation:

$$\mu_1 \sim \mu_2 \iff d\left(\frac{\mu_1}{|\mu_1|}, \frac{\mu_2}{|\mu_2|}\right) < 1$$

Note that  $\frac{\mu_1}{|\mu_1|}$  and  $\frac{\mu_2}{|\mu_2|}$  are distributions. We have that  $\mu_1 \equiv \mu_2$  implies  $\mu_1 \sim \mu_2$ . The relation  $\sim$  is reflexive, symmetric, but in general not transitive. We observe:

**Proposition 18.** Let  $\mu_1 \sim \mu_2$ . Let  $w \in \Sigma^*$  such that  $\text{supp}(\pi_1^w) \supseteq \text{supp}(\mu_1)$  and  $\text{supp}(\pi_2^w) \supseteq \text{supp}(\mu_2)$ . Then  $d(\pi_1, \pi_2) < 1$ .

*Proof.* Since  $\mu_1 \sim \mu_2$ , we have  $d(\rho_1, \rho_2) < 1$  for the distributions  $\rho_1 := \mu_1/|\mu_1|$  and  $\rho_2 := \mu_2/|\mu_2|$ . By Proposition 17 there is a word  $v \in \Sigma^*$  and subdistributions  $\nu_1, \nu_2$  with  $|\nu_1| = |\nu_2| > 0$  and  $\rho_1^v \geq \nu_1 \equiv \nu_2 \leq \rho_2^v$ . Since  $\text{supp}(\pi_i^w) \supseteq \text{supp}(\mu_i) = \text{supp}(\rho_i)$  holds for  $i \in \{1, 2\}$ , we get  $\pi_1^{wv} \geq a\nu_1 \equiv a\nu_2 \leq \pi_2^{wv}$  for some small enough  $a > 0$ . Using Proposition 17 again it follows that  $d(\pi_1, \pi_2) < 1$ .  $\square$



The following proposition states two structural properties of the relation  $\sim$  which can be proved using the fact that  $d(\pi_1, \pi_2) = 1$  implies that there is a “maximizing” event  $E$  with  $\pi_1(E) = 1$  and  $\pi_2(E) = 0$ , see Theorem 11.

**Proposition 19.** *We have the following.*

- (a) Let  $\mu_1 \equiv \mu_2$ . Let  $\nu_1 \leq \mu_1$  with  $|\nu_1| > 0$ . Then  $\nu_1 \sim \mu_2$ .
- (b) Let  $\mu_1 \sim \mu_2$ . Then there is  $q \in \text{supp}(\mu_1)$  with  $\delta_q \sim \mu_2$ .

*Proof.*

- (a) Towards a contradiction suppose that  $d(\nu_1/|\nu_1|, \mu_2/|\mu_2|) = 1$ . Then by Theorem 11 there is an event  $E \subseteq \Sigma^\omega$  with  $\frac{\nu_1}{|\nu_1|}(E) = 0$  and  $\frac{\mu_2}{|\mu_2|}(E) = 1$ , i.e.,  $\nu_1(E) = 0$  and  $\mu_2(E) = |\mu_2|$ . We have:

$$\begin{aligned}
|\mu_2| &= \mu_2(E) \\
&= \mu_1(E) && (\text{as } \mu_1 \equiv \mu_2) \\
&= (\mu_1 - \nu_1)(E) + \nu_1(E) && (\text{as } \nu_1 \leq \mu_1) \\
&= (\mu_1 - \nu_1)(E) && (\text{as } \nu_1(E) = 0) \\
&\leq |\mu_1 - \nu_1| \\
&= |\mu_1| - |\nu_1| && (\text{as } \nu_1 \leq \mu_1) \\
&< |\mu_1| && (\text{as } |\nu_1| > 0) \\
&= |\mu_2| && (\text{as } \mu_1 \equiv \mu_2),
\end{aligned}$$

which is a contradiction. Hence  $\nu_1 \sim \mu_2$ .

- (b) Suppose that for all  $q \in \text{supp}(\mu_1)$  we have  $\delta_q \not\sim \mu_2$ , i.e.,  $d(\delta_q, \mu_2/|\mu_2|) = 1$ . By Theorem 11 for all  $q \in \text{supp}(\mu_1)$  there is an event  $E_q \subseteq \Sigma^\omega$  with  $\delta_q(E_q) = 1$  and  $\mu_2/|\mu_2|(E_q) = 0$ . Consider the event

$$E := \bigcup_{q \in \text{supp}(\mu_1)} E_q.$$

For all  $q \in \text{supp}(\mu_1)$  we have  $\delta_q(E) \geq \delta_q(E_q) = 1$ , so  $\delta_q(E) = 1$ . Hence,

$$\mu_1(E) = \sum_{q \in \text{supp}(\mu_1)} \mu_1(q) \delta_q(E) = \sum_{q \in \text{supp}(\mu_1)} \mu_1(q) = |\mu_1|.$$

On the other hand, by a union bound, we have

$$\mu_2(E) \leq \sum_{q \in \text{supp}(\mu_1)} \mu_2(E_q) = 0.$$

If  $|\mu_2| > 0$ , then by the definition of the distance we have

$$d\left(\frac{\mu_1}{|\mu_1|}, \frac{\mu_2}{|\mu_2|}\right) \geq \frac{\mu_1}{|\mu_1|}(E) - \frac{\mu_2}{|\mu_2|}(E) = 1 - 0 = 1,$$

so  $\mu_1 \not\sim \mu_2$ . If  $|\mu_2| = 0$ , then by the definition of  $\sim$  it also follows  $\mu_1 \not\sim \mu_2$ .  $\square$

For distributions  $\pi_1, \pi_2$  we define a set  $R^{\pi_1, \pi_2} \subseteq Q \times Q$ :

$$R^{\pi_1, \pi_2} := \{(r_1, r_2) \in Q \times Q \mid \exists w \in \Sigma^* : r_1 \in \text{supp}(\pi_1^w) \text{ and } r_2 \in \text{supp}(\pi_2^w)\}$$

This set can be computed in polynomial time:

**Lemma 20.** *Let  $\pi_1, \pi_2$  be distributions. Define a directed graph  $G$  as follows. The vertex set is  $Q \times Q$ . There is an edge from  $(q_1, q_2) \in Q \times Q$  to  $(r_1, r_2) \in Q \times Q$  if there is a  $a \in \Sigma$  with  $M(a)(q_1, r_1) > 0$*

**Algorithm 1** Polynomial-time algorithm for deciding the distance-1 problem.

**procedure** distance1

**input:** LMC  $\mathcal{M} = (Q, \Sigma, M)$

initial distributions  $\pi_1, \pi_2 \in [0, 1]^Q$

**output:**  $d(\pi_1, \pi_2) = 1$  or  $d(\pi_1, \pi_2) < 1$

compute  $\mathcal{B} \subseteq \mathbb{Q}^{2|Q|}$  from Proposition 3 (b)

compute  $R^{\pi_1, \pi_2}$  by graph reachability (Lemma 20)

for  $r_1 \in Q$  do

if there exist subdistributions  $\mu_1, \mu_2$  with

$r_1 \in \text{supp}(\mu_1)$  and  $\text{supp}(\mu_2) \subseteq R_{r_1}^{\pi_1, \pi_2}$   
and  $\forall b \in \mathcal{B} : (\mu_1 \mu_2) \cdot b = 0$

(\* this can be decided using linear programming \*)

then return “ $d(\pi_1, \pi_2) < 1$ ”

fi

od

return “ $d(\pi_1, \pi_2) = 1$ ”

and  $M(a)(q_2, r_2) > 0$ . Then we have:

$$R^{\pi_1, \pi_2}$$

$$= \{(r_1, r_2) \in Q \times Q \mid \exists q_1 \in \text{supp}(\pi_1) \exists q_2 \in \text{supp}(\pi_2) :$$

$$(r_1, r_2) \text{ is reachable from } (q_1, q_2) \text{ in } G\}$$

As a consequence,  $R^{\pi_1, \pi_2}$  can be computed in polynomial time using graph reachability.

The proof of Lemma 20 is straightforward by induction. For  $r_1 \in Q$  we define the projection  $R_{r_1}^{\pi_1, \pi_2} := \{r_2 \in Q \mid (r_1, r_2) \in R^{\pi_1, \pi_2}\}$ . We are ready to show the main theorem of the section:

**Theorem 21.** *Let  $\pi_1, \pi_2$  be distributions. Then  $d(\pi_1, \pi_2) < 1$  holds if and only if there are  $r_1 \in Q$  and subdistributions  $\mu_1, \mu_2$  such that*

$$\mu_1 \equiv \mu_2 \text{ and } r_1 \in \text{supp}(\mu_1) \text{ and } \text{supp}(\mu_2) \subseteq R_{r_1}^{\pi_1, \pi_2}.$$

*Proof. ( $\implies$ )* Let  $d(\pi_1, \pi_2) < 1$ . Then by Proposition 17 there are  $w \in \Sigma^*$  and subdistributions  $\mu_1, \mu_2$  with  $\mu_1 \leq \pi_1^w$  and  $\mu_2 \leq \pi_2^w$  and  $\mu_1 \equiv \mu_2$  and  $|\mu_1| = |\mu_2| > 0$ . By the definition of  $R^{\pi_1, \pi_2}$  we have  $(r_1, r_2) \in R^{\pi_1, \pi_2}$  for all  $r_1 \in \text{supp}(\mu_1)$  and all  $r_2 \in \text{supp}(\mu_2)$ . Choose any  $r_1 \in \text{supp}(\mu_1)$ . Then  $\text{supp}(\mu_2) \subseteq R_{r_1}^{\pi_1, \pi_2}$ .

*( $\impliedby$ )* For the converse, let  $r_1 \in Q$  and  $\mu_1, \mu_2$  be subdistributions such that

$$\mu_1 \equiv \mu_2 \text{ and } r_1 \in \text{supp}(\mu_1) \text{ and } \text{supp}(\mu_2) \subseteq R_{r_1}^{\pi_1, \pi_2}.$$

By Proposition 19 (a) we have  $\delta_{r_1} \sim \mu_2$ . Hence, by Proposition 19 (b) there is  $r_2 \in \text{supp}(\mu_2)$  with  $\delta_{r_1} \sim \delta_{r_2}$ . Since  $(r_1, r_2) \in R^{\pi_1, \pi_2}$ , we have by the definition of  $R^{\pi_1, \pi_2}$  that there is  $w \in \Sigma^*$  with  $r_1 \in \text{supp}(\pi_1^w)$  and  $r_2 \in \text{supp}(\pi_2^w)$ . As  $\delta_{r_1} \sim \delta_{r_2}$ , Proposition 18 implies  $d(\pi_1, \pi_2) < 1$ .  $\square$

We highlight the algorithmic nature of Theorem 21 in Algorithm 1.

## 8. Related Work

Two LMCs have distance 0 if and only if they are language equivalent. We have discussed works on language equivalence in the introduction. The papers [14] and [4] are closest to ours. They investigate the  $L_p$ -distance between two hidden Markov models [14] and two probabilistic automata [4]. Those models are similar to ours. The main difference is that in their models no letters are emitted once a special *end state* is reached, and the transition structure of the chains guarantees that an end state is eventually reached with probability 1. (Our model is more general, as one can make the

LMC emit an infinite sequence of a special “end letter” once the end state is reached.) So those models induce a probability distribution over  $\Sigma^*$ , which makes the sample space countable. As mentioned in the introduction, the  $L_1$ -distance is then twice the total variation distance, so the hardness results from [4, 14] become available: it is NP-hard to “compute” the  $L_1$ - and  $L_\infty$ -distance [14] and the  $L_p$ -distance for odd  $p$  [4], but recall our discussion after Proposition 12 about irrational distances. We note that the example from Figure 3, showing the existence of irrational distances, can be easily framed in their models. It is also shown in [4] that it is NP-hard to approximate the  $L_1$ -distance within an additive error, and that the  $L_p$ -distance can be computed in polynomial time for even  $p$ .

The total variation distance in LMCs is considered in [3], where the authors give an upper bound on the total variation distance in terms of the *bisimilarity pseudometric* defined in [5]. Bisimilarity is a “structural” (i.e., based on the emitted letters and the states) notion of equivalence of LMCs, whereas language equivalence is purely “semantical” (based only on the emitted letters). Accordingly, the bisimilarity pseudometric defines a *branching-time* distance while the total variance distance defines a *linear-time* distance. The authors of [3] prove a quantitative analogue of the fact that bisimilarity implies language equivalence: they prove that the bisimilarity pseudometric, which can be computed in polynomial time [3], is an upper bound on the total variation distance which we discuss here.

## 9. Conclusions and Open Problems

In this paper we have developed a theory of the total variation distance between two LMCs. Two important theoretical results of this paper are summarized as:

- (1) By considering longer and longer prefix words, one can define two sequences  $(1 - \min(i))_{i \in \mathbb{N}}$  and  $(1 - \text{con}(i))_{i \in \mathbb{N}}$  that converge to the distance from below and above, respectively.
- (2) Using the martingale convergence theorem one can show that there is always a maximizing event, and we have explicitly exhibited one.

These results have algorithmic consequences. Our main algorithmic result is a procedure that decides the distance-1 problem in polynomial time. The result (1) also leads to an algorithm for approximating the distance with arbitrary precision. We have also shown that the distance can be irrational, and we have given lower complexity bounds for the threshold-distance problem: it is NP-hard and hard for the square-root-sum problem.

The complexity and even the decidability of the threshold-distance problem are open problems. A theoretical question is whether the distance is always algebraic. We have established a connection to Bernoulli convolutions, the long history of which may hint at the difficulty of solving the mentioned open problems.

## Acknowledgements.

We would like to thank Christoph Haase for pointing us to [10], James Worrell for valuable discussions, and anonymous referees for helpful comments. Stefan Kiefer is supported by a Royal Society University Research Fellowship.

## References

- [1] E. Allender, P. Bürgisser, J. Kjeldgaard-Pedersen, and P. B. Miltersen. On the complexity of numerical analysis. *SIAM J. Comput.*, 38(5): 1987–2006, 2009.
- [2] D. Bailey, J. Borwein, N. Calkin, R. Girgensohn, D. Luke, and V. Moll. *Experimental Mathematics in Action*. Wellesley, 2007.
- [3] D. Chen, F. van Breugel, and J. Worrell. On the complexity of computing probabilistic bisimilarity. In L. Birkedal, editor, *FoSSaCS*, volume 7213 of *Lecture Notes in Computer Science*, pages 437–451. Springer, 2012. ISBN 978-3-642-28728-2.
- [4] C. Cortes, M. Mohri, and A. Rastogi.  $L_p$  distance and equivalence of probabilistic automata. *International Journal of Foundations of Computer Science*, 18(04):761–779, 2007.
- [5] J. Desharnais, V. Gupta, R. Jagadeesan, and P. Panangaden. Metrics for labelled Markov processes. *Theoretical Computer Science*, 318(3): 323–354, 2004.
- [6] L. Doyen, T. Henzinger, and J.-F. Raskin. Equivalence of labeled Markov chains. *International Journal of Foundations of Computer Science*, 19(3):549–563, 2008.
- [7] K. Etessami and M. Yannakakis. On the complexity of nash equilibria and other fixed points. *SIAM J. Comput.*, 39(6):2531–2597, 2010.
- [8] M. Garey, R. Graham, and D. Johnson. Some NP-complete geometric problems. In *STOC*, pages 10–22. ACM, 1976.
- [9] A. Gibbs and F. Su. On choosing and bounding probability metrics. *International Statistical Review*, 70(3):419–435, 2002.
- [10] R. Graham, D. Knuth, and O. Patashnik. *Concrete Mathematics*. Addison-Wesley, second edition, 1989.
- [11] S. Kiefer, A. Murawski, J. Ouaknine, B. Wachter, and J. Worrell. Language equivalence for probabilistic automata. In *CAV*, volume 6806 of *LNCS*, pages 526–540, 2011.
- [12] S. Kiefer, A. Murawski, J. Ouaknine, B. Wachter, and J. Worrell. On the complexity of equivalence and minimisation for Q-weighted automata. *Logical Methods in Computer Science (LMCS)*, 9(1:8):1–22, 2013.
- [13] T. Lengyel. A combinatorial identity and the world series. *SIAM Review*, 35(2):294–297, 1993.
- [14] R. Lyngsø and C. Pedersen. The consensus string problem and the complexity of comparing hidden markov models. *J. Comput. Syst. Sci.*, 65(3):545–569, 2002.
- [15] A. Paz. *Introduction to Probabilistic Automata*. Academic Press, 1971.
- [16] Y. Peres, W. Schlag, and B. Solomyak. Sixty years of Bernoulli convolutions. In *Fractal Geometry and Stochastics II*, volume 46 of *Progress in Probability*, pages 39–65. Birkhäuser Basel, 2000.
- [17] M. O. Rabin. Probabilistic automata. *Information and Control*, 6(3): 230–245, 1963.
- [18] M.-P. Schützenberger. On the definition of a family of automata. *Inf. and Control*, 4:245–270, 1961.
- [19] R. van Glabbeek, S. Smolka, and B. Steffen. Reactive, generative and stratified models of probabilistic processes. *Inf. Comput.*, 121(1):59–80, 1995.
- [20] D. Williams. *Probability with Martingales*. Cambridge University Press, 1991.

## A. Proof of Proposition 3

### Proposition 3.

- (a) We have  $\pi_1 \equiv \pi_2$  if and only if  $d(\pi_1, \pi_2) = 0$ .
- (b) One can compute in polynomial time a set  $\mathcal{B} \subseteq \mathbb{Q}^{2|Q|}$  of column vectors, with  $|\mathcal{B}| \leq 2|Q|$ , such that for all subdistributions  $\mu_1, \mu_2$  we have  $\mu_1 \equiv \mu_2$  if and only if  $(\mu_1 \ \mu_2) \cdot b = 0$  holds for all  $b \in \mathcal{B}$ . Here,  $(\mu_1 \ \mu_2) \in [0, 1]^{2|Q|}$  is the row vector obtained by gluing  $\mu_1, \mu_2$  together. (Note that  $(\mu_1 \ \mu_2) \cdot b$  is a scalar.)
- (c) We have  $\mu_1 \equiv \mu_2$  if and only if  $|\mu_1^w| = |\mu_2^w|$  holds for all  $w \in \Sigma^*$  with  $|w| = 2|Q|$ .
- (d) It is decidable in polynomial time whether  $\mu_1 \equiv \mu_2$  holds. Hence it is also decidable in polynomial time whether  $d(\pi_1, \pi_2) = 0$  holds.

*Proof.*

(a) Immediate from the definitions.

(b) Define

$$\eta := (\underbrace{1, \dots, 1}_{|Q| \text{ times}}, \underbrace{-1, \dots, -1}_{|Q| \text{ times}})^T,$$

where the superscript  $T$  denotes transpose. According to the definitions, we have  $\mu_1 \equiv \mu_2$  if and only if we have

$$(\mu_1 \ \mu_2) \cdot \begin{pmatrix} M(w) & 0 \\ 0 & M(w) \end{pmatrix} \cdot \eta = 0$$

for all  $w \in \Sigma^*$ . We write

$$G := \left\{ \begin{pmatrix} M(w) & 0 \\ 0 & M(w) \end{pmatrix} \cdot \eta \mid w \in \Sigma^* \right\}.$$

Observe that  $G$  is a (column) vector space. As the vectors are  $2|Q|$ -dimensional, a basis of  $G$  contains at most  $2|Q|$  vectors. It is shown, e.g., in [6] that one can compute a basis  $\mathcal{B}$  for  $G$  in  $O(|Q|^3)$  time. It follows that  $\mu_1 \equiv \mu_2$  holds if and only if we have  $(\mu_1 \ \mu_2) \cdot b = 0$  for all  $b \in \mathcal{B}$ .

- (c) The direction “ $\implies$ ” is immediate from the definitions. For the converse, let  $\mu_1 \not\equiv \mu_2$ . By the linear-algebra argument from part (b) there is a word  $u \in \Sigma^*$  with  $|u| \leq 2|Q|$  and  $|\mu_1^u| \neq |\mu_2^u|$ . Towards a contradiction suppose that for all  $v \in \Sigma^*$  with  $|v| = 2|Q| - |u|$  we have  $|\mu_1^{uv}| = |\mu_2^{uv}|$ . Now we have:

$$\begin{aligned} |\mu_1^u| &= \sum_{v \in \Sigma^{2|Q|-|u|}} |\mu_1^{uv}| && \text{(as } \sum_{a \in \Sigma} M(a) \text{ is stochastic)} \\ &= \sum_{v \in \Sigma^{2|Q|-|u|}} |\mu_2^{uv}| && \text{(by assumption)} \\ &= |\mu_2^u| && \text{(as } \sum_{a \in \Sigma} M(a) \text{ is stochastic)} \end{aligned}$$

This is a contradiction. So there is  $v \in \Sigma^*$  with  $|v| = 2|Q| - |u|$  and  $|\mu_1^{uv}| \neq |\mu_2^{uv}|$ .

- (d) Immediate from part (b).

□

## B. Proof of Theorem 7

Recall that for a (random) run  $r \in \Sigma^\omega$  we write  $r_i \in \Sigma^i$  for the length- $i$  prefix of  $r$ . We first prove the following lemma:

**Lemma 22.** *For all  $\varepsilon > 0$  we have*

$$\pi_1(\bar{L} > 0) = \pi_1(\bar{L} > 0 \wedge \exists i \in \mathbb{N} : \min(r_i) \leq (1 + \varepsilon) \text{con}(r_i)).$$

*Proof.* For distributions  $\rho_1, \rho_2$  define  $\tilde{d}(\rho_1, \rho_2) := \max_{w \in \Sigma^{2|Q|}} (|\rho_1^w| - |\rho_2^w|) = \max_{w \in \Sigma^{2|Q|}} (||\rho_1^w| - |\rho_2^w||)$ . Clearly  $\tilde{d}(\rho_1, \rho_2) \geq 0$ .

For a given run  $r$  we define  $\rho_{1,i}$  and  $\rho_{2,i}$  for all  $i \in \mathbb{N}$ : let  $\rho_{1,i} := \pi_1^{r_i} / |\pi_1^{r_i}|$  and  $\rho_{2,i} := \pi_2^{r_i} / |\pi_2^{r_i}|$ . Intuitively,  $\rho_{1,i}$  (resp.  $\rho_{2,i}$ ) is the state distribution in the first (resp. second) LMC, conditioned under having emitted the run prefix  $r_i$ . Define  $\tilde{u}_i \in \Sigma^{2|Q|}$  so that  $\tilde{d}(\rho_{1,i}, \rho_{2,i}) = |\rho_{1,i}^{\tilde{u}_i}| - |\rho_{2,i}^{\tilde{u}_i}|$ . For arbitrary  $i \in \mathbb{N}$  and  $\delta > 0$  define the event  $E_{i,\delta} := \{\tilde{d}(\rho_{1,i}, \rho_{2,i}) \geq \delta\}$ . For any run  $r \in E_{i,\delta}$  we have  $|\rho_{1,i}^{\tilde{u}_i}| - |\rho_{2,i}^{\tilde{u}_i}| \geq \delta$  and hence

$$|\rho_{1,i}^{\tilde{u}_i}| \geq \delta \quad \text{and} \tag{18}$$

$$\frac{|\rho_{2,i}^{\tilde{u}_i}|}{|\rho_{1,i}^{\tilde{u}_i}|} \leq 1 - \frac{\delta}{|\rho_{1,i}^{\tilde{u}_i}|} \leq 1 - \delta. \tag{19}$$

It follows that

$$\begin{aligned}
\delta &\leq \pi_1(r_{i+2|Q|} = r_i \tilde{u}_i \mid E_{i,\delta}) && \text{by (18)} \\
&\leq \pi_1 \left( \begin{array}{c} |\pi_2^{r_{i+2|Q|}}| = |\pi_2^{r_i \tilde{u}_i}| = |\pi_2^{r_i}| |\rho_{2,i}^{\tilde{u}_i}| \\ |\pi_1^{r_{i+2|Q|}}| = |\pi_1^{r_i \tilde{u}_i}| = |\pi_1^{r_i}| |\rho_{1,i}^{\tilde{u}_i}| \end{array} \wedge \mid E_{i,\delta} \right) \\
&\leq \pi_1(L_{i+2|Q|} \leq (1-\delta)L_i \mid E_{i,\delta}) && \text{by (19)}
\end{aligned}$$

In words: for those runs in  $E_{i,\delta}$ , the probability of a decrease of  $L_i$  in the next  $2|Q|$  steps by at least  $\delta L_i$  is bounded below by  $\delta$ . It follows that if  $\tilde{d}(\rho_{1,i}, \rho_{2,i}) \geq \delta$  holds for infinitely many  $i$ , a positive limit  $\bar{L}$  exists with probability 0:

$$\pi_1(\bar{L} > 0 \wedge \tilde{d}(\rho_{1,i}, \rho_{2,i}) \geq \delta \text{ holds for infinitely many } i) = 0.$$

Since  $\delta > 0$  was chosen arbitrarily, we have:

$$\pi_1(\bar{L} > 0 \wedge \lim_{i \rightarrow \infty} \tilde{d}(\rho_{1,i}, \rho_{2,i}) = 0) = \pi_1(\bar{L} > 0). \quad (20)$$

Consider a run with  $\lim_{i \rightarrow \infty} \tilde{d}(\rho_{1,i}, \rho_{2,i}) = 0$ . Since the pairs  $(\rho_{1,i}, \rho_{2,i})$  are elements of the compact set  $C = \{(\rho_1, \rho_2) \mid |\rho_1| = |\rho_2| = 1\}$ , by the Bolzano-Weierstrass theorem, there is a subsequence  $i(0) < i(1) < \dots$  and distributions  $\rho_{1,*}, \rho_{2,*}$  such that

$$\lim_{j \rightarrow \infty} (\rho_{1,i(j)}, \rho_{2,i(j)}) = (\rho_{1,*}, \rho_{2,*}).$$

It follows that for all  $\varepsilon > 0$  there is  $i \in \mathbb{N}$  with

$$\rho_{1,*} \leq (1+\varepsilon)\rho_{1,i} = (1+\varepsilon)\pi_1^{r_i}/|\pi_1^{r_i}| \quad \text{and} \quad \rho_{2,*} \leq (1+\varepsilon)\rho_{2,i} = (1+\varepsilon)\pi_2^{r_i}/|\pi_2^{r_i}|$$

and hence

$$\min(r_i)\rho_{1,*} \leq (1+\varepsilon)\pi_1^{r_i} \quad \text{and} \quad \min(r_i)\rho_{2,*} \leq (1+\varepsilon)\pi_2^{r_i}.$$

Since  $\tilde{d}$  is a continuous function on  $C$ , we have  $\tilde{d}(\rho_{1,*}, \rho_{2,*}) = \tilde{d}(\lim_{i \rightarrow \infty} \rho_{1,i}, \lim_{i \rightarrow \infty} \rho_{2,i}) = \lim_{i \rightarrow \infty} \tilde{d}(\rho_{1,i}, \rho_{2,i}) = 0$ . Hence, by Proposition 3 (c), we have  $\rho_{1,*} \equiv \rho_{2,*}$ , and so  $\min(r_i)\rho_{1,*} \equiv \min(r_i)\rho_{2,*}$ . It follows

$$\min(r_i) \leq (1+\varepsilon)\text{con}(r_i). \quad (21)$$

Using those considerations we obtain:

$$\begin{aligned}
\pi_1(\bar{L} > 0) &= \pi_1(\bar{L} > 0 \wedge \lim_{i \rightarrow \infty} \tilde{d}(\rho_{1,i}, \rho_{2,i}) = 0) && \text{by (20)} \\
&\leq \pi_1(\bar{L} > 0 \wedge \exists i \in \mathbb{N} : \min(r_i) \leq (1+\varepsilon)\text{con}(r_i)) && \text{by (21)} \\
&\leq \pi_1(\bar{L} > 0)
\end{aligned}$$

□

Now we are ready to prove that the limits  $\min(\infty)$  and  $\text{con}(\infty)$  coincide:

**Lemma 23.** *We have  $\min(\infty) = \text{con}(\infty)$ .*

*Proof.* Considering Proposition 4 (c) it suffices to show  $\min(\infty) \leq \text{con}(\infty)$ . Towards a contradiction, suppose this does not hold. Then we have  $\min(\infty) > \text{con}(\infty)$ . So there exists  $\delta \in (0, 1]$  so that for all  $k' \in \mathbb{N}$  we have

$$\min(k') > \text{con}(k') + 4\delta. \quad (22)$$

In the following definition we write  $v \prec w$  to denote that  $v \in \Sigma^*$  is a proper prefix of  $w \in \Sigma^*$ . Define

$$H_\delta := \{w \in \Sigma^* \mid \min(w) \leq (1+\delta)\text{con}(w) \wedge \forall v \prec w : \min(v) > (1+\delta)\text{con}(v)\}.$$

By (22),  $H_\delta \neq \emptyset$ . By Lemma 22 there is  $k_0 \in \mathbb{N}$  so that for all  $k \geq k_0$  we have

$$\pi_1(\bar{L} > 0 \wedge \forall i \leq k : r_i \notin H_\delta) \leq \delta. \quad (23)$$

Using Proposition 6, there is  $k_1 \in \mathbb{N}$  such that for all  $k \geq k_1$  we have

$$\pi_1(\bar{L} = 0 \wedge L_k > \delta) \leq \delta. \quad (24)$$

Choose  $k \geq \max\{k_0, k_1\}$ , so (23) and (24) hold. Define the partition of  $\Sigma^k$  in  $\Sigma^k = W_1 \cup W_2 \cup W_3$  with

$$\begin{aligned}
W_1 &:= \{w \in \Sigma^k \mid |\pi_2^w|/|\pi_1^w| \leq \delta\} \\
W_2 &:= \{w \in \Sigma^k \mid |\pi_2^w|/|\pi_1^w| > \delta \wedge \forall i \leq k : r_i \notin H_\delta\} \\
W_3 &:= \{w \in \Sigma^k \mid |\pi_2^w|/|\pi_1^w| > \delta \wedge \exists i \leq k : r_i \in H_\delta\}
\end{aligned}$$

We have:

$$\begin{aligned}
\sum_{w \in W_1} \min(w) &= \sum_{w \in W_1} |\pi_2^w| \leq \sum_{w \in W_1} \delta |\pi_1^w| \leq \delta \\
\sum_{w \in W_2} \min(w) &\leq \sum_{w \in W_2} |\pi_1^w| = \pi_1(W_2 \Sigma^\omega) \\
&= \pi_1(\bar{L} > 0 \wedge W_2 \Sigma^\omega) \\
&\quad + \pi_1(\bar{L} = 0 \wedge W_2 \Sigma^\omega) \\
&\leq \delta + \delta = 2\delta && \text{(by (23) and (24))} \\
\sum_{w \in W_3} \min(w) &\leq \sum_{w \in H_\delta \text{ s.t. } |w| \leq k} \min(w) && \text{(Prop. 4 (b))} \\
&\leq \sum_{w \in H_\delta \text{ s.t. } |w| \leq k} (1 + \delta) \text{con}(w) && \text{(def. of } H_\delta) \\
&\leq (1 + \delta) \text{con}(k) \leq \text{con}(k) + \delta && \text{(Prop. 4 (b))}
\end{aligned}$$

Adding those inequalities yields

$$\min(k) = \sum_{w \in W_1 \cup W_2 \cup W_3} \min(w) \leq \text{con}(k) + 4\delta,$$

thus contradicting (22), as desired.  $\square$

Now Theorem 7 from the main body of the paper follows:

**Theorem 7.** *We have*

$$1 - \min(\infty) = d(\pi_1, \pi_2) = 1 - \text{con}(\infty).$$

*Proof.* Immediate by combining (4) and Lemma 23.  $\square$

## C. Proofs of Section 6

### C.1 Proof of Proposition 14

We prove Proposition 14 from the main body of the paper:

**Proposition 14.** *The threshold-distance problem is NP-hard with respect to Turing reductions.*

*Proof.* We modify the proof of [14, Section 6]. Let us first sketch some important features of that proof. It is a reduction from the *clique* decision problem: Given a graph  $G = (V, E)$  and a threshold  $t \in \mathbb{N}$ , decide whether  $G$  has a clique of size at least  $t$ . The clique decision problem is known to be NP-complete.

The authors of [14] describe an LMC  $\mathcal{M}_G$ , computed from  $G$ , such that  $\mathcal{M}_G$  emits substrings of the word  $a_1 a_2 \dots a_{|V|}$ , where  $a_1, \dots, a_{|V|} \in \Sigma$ . (In their model, an LMC stops emitting letters once a special end state is reached. This can be simulated in our model by emitting an infinite sequence of a special “end” letter once the end state is reached. In addition, rather than “omitting” letters from  $a_1 a_2 \dots a_{|V|}$  by means of  $\varepsilon$ -transitions, in our model we output a special “blank” symbol. Those changes do not cause problems.) Without loss of generality they assume that  $V = \{1, 2, \dots, |V|\}$ . Define  $\gamma := \sum_{v \in V} 2^{\deg(v)}$ , where  $\deg(v)$  is the number of neighbours of  $v$  in  $G$ . The gadget  $\mathcal{M}_G$  has the following properties for all  $\{i_1, \dots, i_k\} \subseteq \{1, \dots, |V|\}$ :

- If  $i_1 < \dots < i_k$  and  $\{i_1, \dots, i_k\}$  is a clique in  $G$ , then  $a_{i_1} \dots a_{i_k}$  is emitted with probability  $k/\gamma$ ;
- otherwise  $a_{i_1} \dots a_{i_k}$  is emitted with probability 0.

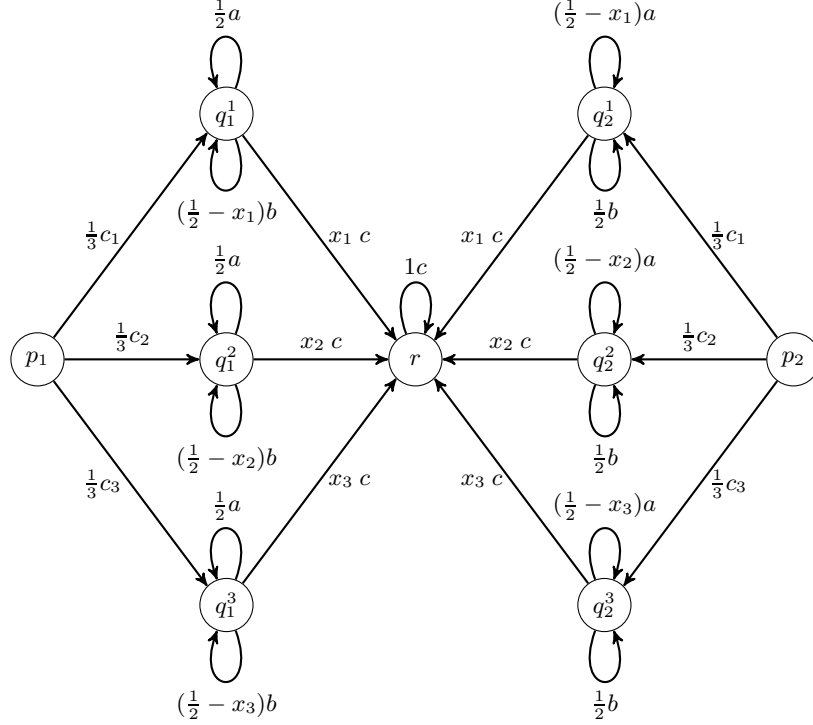
For  $j \in \{0, \dots, |V|\}$ , let  $n_j$  denote the number of cliques of size  $j$ . Note that computing the maximal clique size amounts to computing the maximal  $j$  such that  $n_j \neq 0$ .

Based on the gadget  $\mathcal{M}_G$ , the authors of [14] construct  $|V| + 1$  pairs of LMCs:  $(\mathcal{M}_1^0, \mathcal{M}_2^0), \dots, (\mathcal{M}_1^{|V|}, \mathcal{M}_2^{|V|})$ . For  $i \in \{0, \dots, |V|\}$ , let  $d_i$  denote the distance between  $\mathcal{M}_1^i$  and  $\mathcal{M}_2^i$ . Further, for  $i \leq \gamma/2^{|V|}$  define  $b_i := 1 - \frac{i2^{|V|}}{\gamma}$  and  $c_i := 1/\gamma$ ; for  $i > \gamma/2^{|V|}$  define  $b_i := 1 - \frac{\gamma}{i2^{|V|}}$  and  $c_i := \frac{1}{i2^{|V|}}$ . The pairs  $(\mathcal{M}_1^i, \mathcal{M}_2^i)$  are constructed such that

$$d_i = b_i + c_i \sum_{j=0}^{|V|} n_j |i - j| \quad (25)$$

holds for all  $i \in \{0, \dots, |V|\}$ . It is argued in [14] that once  $d_0, \dots, d_{|V|}$  are known, one can compute  $n_0, \dots, n_{|V|}$  by solving the linear equation system (25). The largest clique size is then the maximal  $j$  such that  $n_j \neq 0$ . This establishes a Turing reduction from the clique decision problem to computing the distance. (Note that by (25) the distances are rational in this reduction.)

For a Turing reduction from the clique decision problem to the threshold-distance problem, it now suffices to argue that one can compute all  $d_i$  with a polynomial number of threshold queries (“is the distance at least  $\tau$ ?”). Indeed, consider again (25). The term  $\sum_{j=0}^{|V|} n_j |i - j|$  is an



**Figure 6.** This LMC is obtained by combining the chain from Figure 3 in parallel  $n = 3$  times. We have  $d(\delta_{p_1}, \delta_{p_2}) = \frac{1}{3} (d(\delta_{q_1^1}, \delta_{q_2^1}) + d(\delta_{q_1^2}, \delta_{q_2^2}) + d(\delta_{q_1^3}, \delta_{q_2^3}))$ .

integer between 0 and  $(|V|+1) \cdot 2^{|V|} \cdot |V|$ . Since  $b_i$  and  $c_i$  are fixed and known for each  $i$ , each  $d_i$  takes one of at most  $(|V|+1) \cdot 2^{|V|} \cdot |V|+1$  values. Using binary search,  $\log_2((|V|+1) \cdot 2^{|V|} \cdot |V|+1)$  distance queries suffice to determine  $d_i$ . Doing this for each  $i$  results in a polynomial number of threshold queries.  $\square$

## C.2 Proof of Theorem 15

We prove Theorem 15 from the main body of the paper.

**Theorem 15.** *There is a polynomial-time many-one reduction from the square-root-sum problem to the threshold-distance problem.*

*Proof.* The construction is by taking the LMC from Figure 3 as a gadget, and joining  $n$  instances of it in parallel. This is sketched for  $n = 3$  in Figure 6. In general we have  $\Sigma = \{c_1, \dots, c_n, a, b, c\}$  and  $Q = \{p_1, p_2, q_1^1, \dots, q_1^n, q_2^1, \dots, q_2^n, r\}$ . Using this construction we have

$$d(\delta_{p_1}, \delta_{p_2}) = \frac{1}{n} \sum_{i=1}^n d(\delta_{q_1^i}, \delta_{q_2^i}). \quad (26)$$

To see this, consider the event

$$E_{\geq} := \{wccc \dots \mid w \in \{a, b\}^*, \#_a(w) \geq \#_b(w)\}$$

from the proof of Proposition 12. From that proof we know

$$d(\delta_{q_1^i}, \delta_{q_2^i}) = \delta_{q_1^i}(E_{\geq}) - \delta_{q_2^i}(E_{\geq}). \quad (27)$$

It follows that we have

$$d(\delta_{p_1}, \delta_{p_2}) = \delta_{p_1}(E) - \delta_{p_2}(E) \quad (28)$$

for

$$E := \{\sigma wccc \dots \mid \sigma \in \{c_1, \dots, c_n\}, w \in \{a, b\}^*, \#_a(w) \geq \#_b(w)\}.$$

From the definition of  $E$  we have

$$\delta_{p_1}(E) = \frac{1}{n} \sum_{i=1}^n \delta_{q_1^i}(E_{\geq}) \quad \text{and} \quad \delta_{p_2}(E) = \frac{1}{n} \sum_{i=1}^n \delta_{q_2^i}(E_{\geq}). \quad (29)$$

It follows:

$$d(\delta_{p_1}, \delta_{p_2}) = \delta_{p_1}(E) - \delta_{p_2}(E) \quad \text{by (28)}$$

$$= \frac{1}{n} \sum_{i=1}^n (\delta_{q_1^i}(E_{\geq}) - \delta_{q_2^i}(E_{\geq})) \quad \text{by (29)}$$

$$= \frac{1}{n} \sum_{i=1}^n d(\delta_{q_1^i}, \delta_{q_2^i}) \quad \text{by (27),}$$

hence (26) is proved.

Recall that the input of the square-root-sum problem is a list of integers  $s_1, \dots, s_n \in \mathbb{N}$  and  $t \in \mathbb{N}$ . Without loss of generality, we can assume that  $s_1, \dots, s_n, t \geq 1$ . The reduction is as follows. Define  $h := 3 \max_{i \in \{1, \dots, n\}} s_i$ . Construct the LMC from Figure 6 with  $x_i := 2s_i/h^2$ . Then we have  $x_i \in (0, 1/2)$  and

$$\frac{1}{2} \sqrt{2x_i} = \frac{1}{h} \sqrt{s_i}. \quad (30)$$

Set the threshold  $\tau := \frac{1}{nh} \cdot t$ . We have:

$$\begin{aligned} d(\delta_{p_1}, \delta_{p_2}) &= \frac{1}{n} \sum_{i=1}^n d(\delta_{q_1^i}, \delta_{q_2^i}) && \text{(by (26))} \\ &= \frac{1}{n} \sum_{i=1}^n \frac{1}{2} \sqrt{2x_i} && \text{(Proposition 12)} \\ &= \frac{1}{nh} \sum_{i=1}^n \sqrt{s_i} && \text{(by (30)).} \end{aligned}$$

It follows that we have  $d(\delta_{p_1}, \delta_{p_2}) \geq \tau$  if and only if  $\sum_{i=1}^n \sqrt{s_i} \geq t$ .  $\square$

## D. Proof of Proposition 16

Recall that by Theorem 7 we have  $d_\theta(x) = 1 - \text{con}(\infty)$ . Let  $\pi_1, \pi_2$  be the initial distributions concentrated on  $p_1, p_2$ , respectively. It is easy to see that for any word  $w \in \Sigma^*$  and any  $\mu_1, \mu_2$  with  $\mu_1 \leq \pi_1^w$  and  $\mu_2 \leq \pi_2^w$  and  $\mu_1 \equiv \mu_2$  we have  $\mu_1(q_1) = \mu_2(q_1)$  and  $\mu_1(q_2) = \mu_2(q_2)$  and  $\mu_1(s) = \mu_2(s) = 0$  for  $s \in \{p_1, p_2, r_1, r_2\}$ . Therefore, writing  $\mu_1 \wedge \mu_2$  for the componentwise minimum of row vectors  $\mu_1, \mu_2$ , we have  $\text{con}(w) = |\pi_1^w \wedge \pi_2^w|$ . Hence  $\text{con}(bw) = 0$  holds for all  $w \in \Sigma^*$ . So we have for all  $k \geq 1$ :

$$\text{con}(k) = \sum_{w \in \Sigma^k} \text{con}(w) = \sum_{w \in \Sigma^{k-1}} \text{con}(aw) = \sum_{w \in \Sigma^{k-1}} |\pi_1^{aw} \wedge \pi_2^{aw}| \quad (31)$$

Note that the states  $q_1, q_2$  cannot be left after they have been entered. This motivates the definition of transition matrices restricted to  $q_1, q_2$ :

$$A = \begin{pmatrix} \frac{1}{2} & \frac{1}{2} - \frac{1}{2\theta} \\ 0 & \frac{1}{2\theta} \end{pmatrix} \quad B = \begin{pmatrix} \frac{1}{2} - \frac{1}{2\theta} & 0 \\ \frac{1}{2} & \frac{1}{2\theta} \end{pmatrix}$$

Define for each  $k \geq 1$  a function  $\text{co}_k : [0, \infty)^2 \times [0, \infty)^2 \rightarrow [0, \infty)$  by

$$\begin{aligned} \text{co}_1(u, v) &= |u \wedge v| \\ \text{co}_{k+1}(u, v) &= \text{co}_k(uA, vA) + \text{co}_k(uB, vB) \end{aligned}$$

By the definition of  $\pi_1^{aw}, \pi_2^{aw}$  and by (31) we have

$$\text{con}(k) = \text{co}_k \left( \left( \frac{1}{2} - x, 0 \right), \left( 0, \frac{1}{2} + x \right) \right) \quad \text{for all } k \geq 1. \quad (32)$$

Define for each  $k \geq 1$  a function  $li_k : \mathbb{R}^2 \rightarrow [0, \infty)$  by

$$\begin{aligned} li_1(z) &= |z| \\ li_{k+1}(z) &= \begin{cases} |z| & \text{if } z \leq (0, 0) \text{ or } z \geq (0, 0) \\ li_k(zA) + li_k(zB) & \text{otherwise,} \end{cases} \end{aligned}$$

where, for  $z_1, z_2 \in \mathbb{R}$ , we write  $|(z_1, z_2)| = |z_1| + |z_2|$  for the  $L_1$ -norm. It follows immediately from the definition that the function  $li_k$  is “almost linear” in the sense that we have

$$li_k(az) = ali_k(z) = li_k(-az) \quad \text{for } a \in [0, \infty) \text{ and } k \in \{1, 2, \dots\}. \quad (33)$$

The following lemma connects  $\text{co}_k$  and  $li_k$ :

**Lemma 24.** *Let  $u \geq (0, 0)$  and  $v \geq (0, 0)$ . Then for all  $k \geq 1$ :*

$$2\text{co}_k(u, v) + li_k(u - v) = |u + v|.$$

*Proof.* We write  $u = (u_1, u_2)$  and  $v = (v_1, v_2)$ . We proceed by induction on  $k$ . For the induction base let  $k = 1$ . Let  $u_1 \leq v_1$  and  $u_2 \leq v_2$ . Then we have  $co_1(u, v) = u_1 + u_2$  and  $li_1(u - v) = v_1 - u_1 + v_2 - u_2$ . Hence  $2co_1(u, v) + li_1(u - v) = u_1 + v_1 + u_2 + v_2 = |u + v|$ . The case  $u_1 \geq v_1$  and  $u_2 \geq v_2$  is similar.

Now let  $u_1 \leq v_1$  and  $u_2 \geq v_2$ . Then we have  $co_1(u, v) = u_1 + v_2$  and  $li_1(u - v) = v_1 - u_1 + u_2 - v_2$ . Hence  $2co_1(u, v) + li_1(u - v) = u_1 + v_1 + u_2 + v_2 = |u + v|$ . The case  $u_1 \geq v_1$  and  $u_2 \leq v_2$  is similar.

For the induction step let  $k \geq 1$ . Then we have:

$$li_{k+1}(u - v) = li_k((u - v)A) + li_k((u - v)B) . \quad (34)$$

Indeed, if  $u \leq v$ , we have

$$\begin{aligned} li_{k+1}(u - v) &= |u - v| && \text{(definition of } li_{k+1}) \\ &= |(u - v)A + (u - v)B| && ((v - u) \geq 0 \text{ and } A + B \text{ is stochastic}) \\ &= |(u - v)A| + |(u - v)B| && ((v - u) \text{ and } A \text{ and } B \text{ are nonnegative}) \\ &= li_k((u - v)A) + li_k((u - v)B) && \text{(by definition of } li_k) \end{aligned}$$

The same equalities hold in the case  $u \geq v$ . If neither  $u \leq v$  nor  $u \geq v$  holds, then we have  $li_{k+1}(u - v) = li_k((u - v)A) + li_k((u - v)B)$  by the definition of  $li_{k+1}$ .

So we have:

$$\begin{aligned} 2co_{k+1}(u, v) + li_{k+1}(u - v) &= 2co_k(uA, vA) + 2co_k(uB, vB) + li_{k+1}(u - v) && \text{(definition of } co_{k+1}) \\ &= 2co_k(uA, vA) + 2co_k(uB, vB) + li_k((u - v)A) + li_k((u - v)B) && \text{(by (34))} \\ &= |uA + vA| + |uB + vB| && \text{(induction hypothesis)} \\ &= |(u + v)(A + B)| && (u, v, A, B \text{ are nonnegative}) \\ &= |u + v| && (A + B \text{ stochastic}) \end{aligned}$$

□

Summarizing the previous development we obtain:

**Lemma 25.** *We have:*

$$d_\theta(x) = \frac{1}{2} + \frac{1}{2} \lim_{k \rightarrow \infty} li_k \left( \left( x - \frac{1}{2}, x + \frac{1}{2} \right) \right)$$

*Proof.*

$$\begin{aligned} d_\theta(x) &= 1 - con(\infty) && \text{(Theorem 7)} \\ &= 1 - \lim_{k \rightarrow \infty} co_k \left( \left( \frac{1}{2} - x, 0 \right), \left( 0, \frac{1}{2} + x \right) \right) && \text{(by (32))} \\ &= \frac{1}{2} + \frac{1}{2} \lim_{k \rightarrow \infty} li_k \left( \left( \frac{1}{2} - x, -\frac{1}{2} - x \right) \right) && \text{(Lemma 24)} \\ &= \frac{1}{2} + \frac{1}{2} \lim_{k \rightarrow \infty} li_k \left( \left( x - \frac{1}{2}, x + \frac{1}{2} \right) \right) && \text{(by (33))} \end{aligned}$$

□

Lemma 25 suggests the definition of a function  $f^{(k)} : \mathbb{R} \rightarrow [0, \infty)$ , for each  $k \geq 1$ , such that

$$f^{(k)}(x) = li_k \left( x - \frac{1}{2}, x + \frac{1}{2} \right) .$$

The following lemma characterizes  $f^{(k)}$  recursively:

**Lemma 26.** *We have  $f^{(1)}(x) = 2|x|$ . Further, for all  $k \geq 1$ :*

$$f^{(k+1)}(x) = \begin{cases} 2|x| & |x| \geq \frac{1}{2} \\ \frac{1}{2\theta} f^{(k)}(\theta x - (\frac{1}{2}\theta - \frac{1}{2})) + \frac{1}{2\theta} f^{(k)}(\theta x + (\frac{1}{2}\theta - \frac{1}{2})) & |x| \leq \frac{1}{2} \end{cases}$$

*Proof.* The equalities  $f^{(1)}(x) = 2|x|$  and  $f^{(k+1)}(x) = 2|x|$  for  $|x| \geq \frac{1}{2}$  follow from the definitions. Further we have:

$$\begin{aligned} \left( x - \frac{1}{2}, x + \frac{1}{2} \right) A &= \frac{1}{2\theta} \left( \theta x - \frac{1}{2}\theta, \theta x - \frac{1}{2}\theta + 1 \right) && \text{and} \\ \left( x - \frac{1}{2}, x + \frac{1}{2} \right) B &= \frac{1}{2\theta} \left( \theta x + \frac{1}{2}\theta - 1, \theta x + \frac{1}{2}\theta + 1 \right) \end{aligned}$$



So it follows for  $|x| \leq \frac{1}{2}$ :

$$\begin{aligned}
f^{(k+1)}(x) &= li_{k+1} \left( x - \frac{1}{2}, x + \frac{1}{2} \right) && \text{(definition of } f^{(k+1)}) \\
&= li_k \left( \frac{1}{2\theta} \left( \theta x - \frac{1}{2}\theta, \theta x - \frac{1}{2}\theta + 1 \right) \right) && \text{(definition of } li_{k+1}) \\
&\quad + li_k \left( \frac{1}{2\theta} \left( \theta x + \frac{1}{2}\theta - 1, \theta x + \frac{1}{2}\theta + 1 \right) \right) \\
&= \frac{1}{2\theta} li_k \left( \theta x - \frac{1}{2}\theta, \theta x - \frac{1}{2}\theta + 1 \right) && \text{(by (33))} \\
&\quad + \frac{1}{2\theta} li_k \left( \theta x + \frac{1}{2}\theta - 1, \theta x + \frac{1}{2}\theta + 1 \right) \\
&= \frac{1}{2\theta} f^{(k)} \left( \theta x - \left( \frac{1}{2}\theta - \frac{1}{2} \right) \right) && \text{(definition of } f^{(k)}) \\
&\quad + \frac{1}{2\theta} f^{(k)} \left( \theta x + \left( \frac{1}{2}\theta - \frac{1}{2} \right) \right)
\end{aligned}$$

□

Now we prove Proposition 16 from the main body of the paper.

**Proposition 16.** For all  $\theta > 1$  we have  $d_\theta(x) = \frac{1}{2} + \frac{1}{2}f_\theta(x)$  for the unique function  $f_\theta : \mathbb{R} \rightarrow \mathbb{R}$  with

$$f_\theta(x) = \begin{cases} -2x & x \leq -\frac{1}{2} \\ \frac{1}{2\theta}f_\theta(\theta x - (\frac{1}{2}\theta - \frac{1}{2})) & x \in [-\frac{1}{2}, \frac{1}{2}] \\ + \frac{1}{2\theta}f_\theta(\theta x + (\frac{1}{2}\theta - \frac{1}{2})) & \\ 2x & x \geq +\frac{1}{2} . \end{cases}$$

*Proof.* We use the Banach fixed-point theorem. Define a complete metric space  $(F, \Delta)$  by

$$F := \{f : \mathbb{R} \rightarrow [0, \infty) \mid f \text{ is continuous and } f(x) = 2|x| \text{ for } |x| \geq \frac{1}{2}\}$$

and the distance metric  $\Delta : F \times F \rightarrow [0, \infty)$  with

$$\Delta(f_1, f_2) := \sup_{|x| \leq \frac{1}{2}} |f_1(x) - f_2(x)| .$$

Fix  $\theta > 1$ . Define the function  $S_\theta : F \rightarrow F$  with

$$S_\theta(f)(x) := \begin{cases} 2|x| & |x| \geq \frac{1}{2} \\ \frac{1}{2\theta}f(\theta x - (\frac{1}{2}\theta - \frac{1}{2})) + \frac{1}{2\theta}f(\theta x + (\frac{1}{2}\theta - \frac{1}{2})) & |x| \leq \frac{1}{2} . \end{cases}$$

We have for all  $f_1, f_2 \in F$  and  $|x| \leq \frac{1}{2}$ :

$$\begin{aligned}
&|S_\theta(f_1)(x) - S_\theta(f_2)(x)| \\
&\leq \frac{1}{2\theta} \left| f_1 \left( \theta x - \left( \frac{1}{2}\theta - \frac{1}{2} \right) \right) - f_2 \left( \theta x - \left( \frac{1}{2}\theta - \frac{1}{2} \right) \right) \right| \\
&\quad + \frac{1}{2\theta} \left| f_1 \left( \theta x + \left( \frac{1}{2}\theta - \frac{1}{2} \right) \right) - f_2 \left( \theta x + \left( \frac{1}{2}\theta - \frac{1}{2} \right) \right) \right| \\
&\leq \frac{1}{2\theta} \Delta(f_1, f_2) + \frac{1}{2\theta} \Delta(f_1, f_2) \\
&= \frac{1}{\theta} \Delta(f_1, f_2)
\end{aligned}$$

It follows that we have  $\Delta(S_\theta(f_1), S_\theta(f_2)) \leq \frac{1}{\theta} \Delta(f_1, f_2)$ , so  $S_\theta$  is contraction mapping. Using the Banach fixed-point theorem and Lemma 26 we obtain that the function sequence  $f^{(1)}, f^{(2)}, \dots$  converges to the (unique) fixed point of  $S_\theta$ , i.e., to the function  $f_\theta$  from the statement of this proposition. It follows:

$$\begin{aligned}
d_\theta(x) &= \frac{1}{2} + \frac{1}{2} \lim_{k \rightarrow \infty} li_k \left( \left( x - \frac{1}{2}, x + \frac{1}{2} \right) \right) && \text{(Lemma 25)} \\
&= \frac{1}{2} + \frac{1}{2} \lim_{k \rightarrow \infty} f^{(k)}(x) && \text{(definition of } f^{(k)}) \\
&= \frac{1}{2} + \frac{1}{2} f_\theta(x) && \text{(as argued above)}
\end{aligned}$$

□

## E. Proof of Lemma 20

We prove Lemma 20 from the main body of the paper.

**Lemma 20.** *Let  $\pi_1, \pi_2$  be distributions. Define a directed graph  $G$  as follows. The vertex set is  $Q \times Q$ . There is an edge from  $(q_1, q_2) \in Q \times Q$  to  $(r_1, r_2) \in Q \times Q$  if there is  $a \in \Sigma$  with  $M(a)(q_1, r_1) > 0$  and  $M(a)(q_2, r_2) > 0$ . Then we have:*

$$\begin{aligned} R^{\pi_1, \pi_2} &= \{(r_1, r_2) \in Q \times Q \mid \exists q_1 \in \text{supp}(\pi_1) \exists q_2 \in \text{supp}(\pi_2) : \\ &\quad (r_1, r_2) \text{ is reachable from } (q_1, q_2) \text{ in } G\} \end{aligned}$$

As a consequence,  $R^{\pi_1, \pi_2}$  can be computed in polynomial time using graph reachability.

*Proof.* Define

$$\begin{aligned} S^{\pi_1, \pi_2} &:= \{(r_1, r_2) \in Q \times Q \mid \exists q_1 \in \text{supp}(\pi_1) \exists q_2 \in \text{supp}(\pi_2) : \\ &\quad (r_1, r_2) \text{ is reachable from } (q_1, q_2) \text{ in } G\}. \end{aligned}$$

We need to show  $S^{\pi_1, \pi_2} = R^{\pi_1, \pi_2}$ .

First we show  $S^{\pi_1, \pi_2} \subseteq R^{\pi_1, \pi_2}$ . For  $k \in \mathbb{N}$  let

$$\begin{aligned} S_k^{\pi_1, \pi_2} &:= \{(r_1, r_2) \in Q \times Q \mid \exists q_1 \in \text{supp}(\pi_1) \exists q_2 \in \text{supp}(\pi_2) : \\ &\quad (r_1, r_2) \text{ is reachable from } (q_1, q_2) \text{ in } G \text{ in } k \text{ steps.}\} \end{aligned}$$

We show by induction on  $k$  that

$$S_k^{\pi_1, \pi_2} \subseteq R^{\pi_1, \pi_2} \quad \text{for all } k \in \mathbb{N}. \quad (35)$$

The case  $k = 0$  is trivial. Let  $k \geq 0$ . Let  $(r'_1, r'_2) \in S_{k+1}^{\pi_1, \pi_2}$ . Then there are  $q_1 \in \text{supp}(\pi_1)$  and  $q_2 \in \text{supp}(\pi_2)$  and  $r_1, r_2 \in Q$  so that there is a path

$$(q_1, q_2) \rightarrow \cdots \rightarrow (r_1, r_2) \rightarrow (r'_1, r'_2)$$

of length  $k + 1$  in  $G$ . By the induction hypothesis there is  $w \in \Sigma^*$  such that  $r_1 \in \text{supp}(\pi_1^w)$  and  $r_2 \in \text{supp}(\pi_2^w)$ . By the definition of  $G$  the presence of an edge from  $(r_1, r_2)$  to  $(r'_1, r'_2)$  implies that there is  $a \in \Sigma$  with  $M(a)(r_1, r'_1) > 0$  and  $M(a)(r_2, r'_2) > 0$ . It follows that we have  $r'_1 \in \text{supp}(\pi_1^{wa})$  and  $r'_2 \in \text{supp}(\pi_2^{wa})$ . Hence  $(r'_1, r'_2) \in R^{\pi_1, \pi_2}$ . This proves (35). Since  $\bigcup_{k \geq 0} S_k^{\pi_1, \pi_2} = S^{\pi_1, \pi_2}$ , we have also shown  $S^{\pi_1, \pi_2} \subseteq R^{\pi_1, \pi_2}$ .

Now we show  $R^{\pi_1, \pi_2} \subseteq S^{\pi_1, \pi_2}$ . For  $w \in \Sigma^*$  let

$$R_w^{\pi_1, \pi_2} := \{(r_1, r_2) \in Q \times Q \mid r_1 \in \text{supp}(\pi_1^w) \text{ and } r_2 \in \text{supp}(\pi_2^w)\}.$$

We show by induction on the length of  $w$  that

$$R_w^{\pi_1, \pi_2} \subseteq S^{\pi_1, \pi_2} \quad \text{for all } w \in \Sigma^*. \quad (36)$$

The case  $|w| = 0$  is trivial. Let  $w \in \Sigma^*$  and  $a \in \Sigma$ . Let  $(r'_1, r'_2) \in R_{wa}^{\pi_1, \pi_2}$ , i.e.,  $r'_i \in \text{supp}(\pi_i^w M(a))$  for  $i \in \{1, 2\}$ . Hence there are  $r_1, r_2 \in Q$  with  $r_i \in \text{supp}(\pi_i^w)$  and  $M(a)(r_i, r'_i) > 0$  for  $i \in \{1, 2\}$ . By the induction hypotheses we have that  $(r_1, r_2)$  is reachable from  $(q_1, q_2)$  in  $G$ . By the definition of  $G$  there is an edge from  $(r_1, r_2)$  to  $(r'_1, r'_2)$  in  $G$ . Hence  $(r'_1, r'_2)$  is reachable from  $(q_1, q_2)$  in  $G$ , so  $(r'_1, r'_2) \in S^{\pi_1, \pi_2}$  and (36) is proved. Since  $\bigcup_{w \in \Sigma^*} R_w^{\pi_1, \pi_2} = R^{\pi_1, \pi_2}$ , we have also shown  $R^{\pi_1, \pi_2} \subseteq S^{\pi_1, \pi_2}$ .  $\square$